

# Hardware Acceleration for HPS Algorithms in Two and Three Dimensions

Owen Melia<sup>1</sup>, Daniel Fortunato<sup>1,2</sup>, Jeremy Hoskins<sup>3</sup>, and Rebecca Willett<sup>3,4,5</sup>

<sup>1</sup>*Center for Computational Mathematics, Flatiron Institute, USA*

<sup>2</sup>*Center for Computational Biology, Flatiron Institute, USA*

<sup>3</sup>*Computational and Applied Mathematics, Department of Statistics, University of Chicago, USA*

<sup>4</sup>*Department of Computer Science, University of Chicago, USA*

<sup>5</sup>*Data Science Institute, University of Chicago, USA*

## Abstract

We provide a flexible, open-source framework for hardware acceleration, namely massively-parallel execution on general-purpose graphics processing units (GPUs), applied to the hierarchical Poincaré–Steklov (HPS) family of algorithms for building fast direct solvers for linear elliptic partial differential equations. To take full advantage of the power of hardware acceleration, we propose two variants of HPS algorithms to improve performance on two- and three-dimensional problems. In the two-dimensional setting, we introduce a novel recomputation strategy that minimizes costly data transfers to and from the GPU; in three dimensions, we modify and extend the adaptive discretization technique of Geldermans and Gillman (2019) to greatly reduce peak memory usage. We provide an open-source implementation of these methods written in JAX, a high-level accelerated linear algebra package, which allows for the first integration of a high-order fast direct solver with automatic differentiation tools. We conclude with extensive numerical examples showing our methods are fast and accurate on two- and three-dimensional problems.

## 1 Introduction

Many problems in scientific computing require solving systems of linear, elliptic partial differential equations (PDEs). Such PDEs can accurately model a variety of physics, such as wave propagation, electrostatics, and diffusion phenomena. Because analytical solutions of these equations are often unknown, the task of computing numerical solutions has been an area of active research for hundreds of years. Today, there are a myriad of numerical solution methods available, and many are tailored to particular classes of equations or to specific use cases. We are most interested in designing methods for settings such as inverse or control problems, where the PDE implicitly defines some functional which, along with its gradient, is evaluated sequentially hundreds or thousands of times in the inner loop of an iterative algorithm.

In these settings, fast direct solvers (Martinsson, 2019) are a compelling choice. These solvers are able to rapidly compute a high-accuracy solution operator and can rapidly evaluate the solution given new data by applying the solution operator, often at the cost of a few matrix-vector multiplications. Fast direct solvers are also preferable for certain PDEs with oscillatory solutions (Gillman et al., 2015), especially ones modeling wave propagation, as they do not incur a data-dependent iteration complexity cost associated with iterative solvers, which can be quite large (Ernst and Gander, 2012).

Recently, scientific computing has undergone a paradigm shift with the advent of general-purpose hardware accelerators, such as GPUs. These hardware accelerators allow for massively

parallel computation—they have thousands of processor cores on a single chip—but have strict memory constraints, a resource profile very different from standard multicore CPU architectures. This paper explores hardware acceleration of fast direct solvers and introduces new methods to facilitate this acceleration.

In particular, we focus on the hierarchical Poincaré–Steklov family of algorithms (Martinsson, 2013; Gillman and Martinsson, 2014; Gillman et al., 2015), a class of direct solution methods for variable-coefficient elliptic PDEs. These methods are characterized by a nested dissection approach combined with a high-order composite spectral discretization. We identify the algorithmic structure of these algorithms which makes them amenable to GPU acceleration and introduce new techniques for reducing the memory footprint of these methods. For two-dimensional problems, we introduce a novel recomputation strategy that minimizes data transfer between the GPU and host memory. In three dimensions, we use an adaptive discretization method, which greatly reduces the algorithm’s peak memory complexity. Our numerical examples show these ideas are useful in challenging applied settings such as wave propagation, inverse problems, and molecular biology simulations.

We focus on solving linear, elliptic partial differential equations of the form

$$\mathcal{L}u(x) = f(x), \quad x \in \Omega, \quad (1)$$

$$u(x) = g(x), \quad x \in \partial\Omega. \quad (2)$$

In Equation (1),  $\mathcal{L}$  is a linear, elliptic, second-order partial differential operator with spatially-varying coefficient functions, and  $\Omega$  is a square  $\subset \mathbb{R}^2$  or cube  $\subset \mathbb{R}^3$ . Equation (2) specifies Dirichlet boundary data, but our methods can also solve problems with Robin or Neumann boundary data. In these problems, we assume we can evaluate the differential operator  $\mathcal{L}$ , the source  $f$ , and the boundary data  $g$  at a set of discretization points of our choosing. We represent the solution  $u$  by its restriction to the same set of discretization points and rely on high-order polynomial interpolation to evaluate  $u$  away from the discretization points. In this paper, we refer to vectors with bold lowercase symbols such as  $\mathbf{f}$  and matrices with bold uppercase symbols such as  $\mathbf{A}$ . We use  $x$  for the spatial variable, and when we want to indicate Cartesian coordinates, we use  $(x_1, x_2) \in \mathbb{R}^2$  and  $(x_1, x_2, x_3) \in \mathbb{R}^3$ . We use the subscript  $u_n$  to denote the outward-pointing boundary normal derivative of a function, and we use  $\Delta$  to denote the Laplace operator, the sum of second derivatives in each dimension.

## 1.1 Paper outline and contributions

In Section 2, we discuss related work, including algorithmic development for fast direct solvers and GPU-specific optimizations. In Section 3, we give an overview of the hierarchical Poincaré–Steklov method and discuss the potential for massively parallel implementations of the algorithm. In the rest of the paper, we make the following contributions:

- We optimize data transfer patterns to accelerate our method applied to two-dimensional problems (Section 4.1).
- To alleviate peak memory requirements in three-dimensional problems, we extend the two-dimensional adaptive method of Geldermans and Gillman (2019) to three dimensions, develop the first adaptive 3D GPU-compatible HPS implementation, and provide numerical examples to demonstrate memory and accuracy tradeoffs (Section 4.2).
- We provide a range of numerical examples illustrating the application of our method, focusing on two settings: high-wavenumber scattering problems and the linearized Poisson–Boltzmann equation (Sections 5 and 6).

- We show our proposed algorithm and implementation can be combined easily with standard automatic differentiation software, which makes it particularly amenable to application in optimization, inverse problems, and machine learning contexts (Section 5).
- We make our JAX-based implementation publicly available at <https://github.com/meliao/jaxhps>.

## 2 Related work

HPS algorithms are built on two conceptual building blocks: composite high-order spectral collocation methods (Kopriva, 1998; Pfeiffer et al., 2003; Yang and Hesthaven, 2000), and nested dissection of the computational domain (George, 1973). Composite spectral collocation methods are those which separate the computational domain into a set of disjoint elements, and use a high-order spectral collocation scheme to represent the problem and solution separately on each element. Nested dissection methods break the original problem into a series of subproblems defined on a hierarchy of subdomains. The careful ordering of subproblems reduces the overall computational complexity by leveraging knowledge about properties of the solution, i.e. continuity of the solution and its derivative. Composite spectral collocation and hierarchical matrix decomposition ideas were combined in integral equation methods (Ho and Greengard, 2012) for constant-coefficient PDEs.

Martinsson (2013) first proposed combining these elements in a fast direct solver for variable-coefficient linear elliptic PDEs. The proposed scheme discretizes and merges Dirichlet-to-Neumann (DtN) operators. Gillman and Martinsson (2014) proposed a compression scheme that leverages the structure of these DtN operators to build a solver with  $O(n)$  computational complexity for  $n$  elements. To alleviate the instabilities observed when merging DtN operators for Helmholtz problems, Gillman et al. (2015) proposed a scheme that merges impedance-to-impedance (ItI) operators instead. Further analysis for this scheme was provided in Beck et al. (2022). Modifications for three dimensions have been proposed, including Lucero Lorca et al. (2024); Hao and Martinsson (2016); Kump et al. (2025), which all build solvers for three-dimensional Helmholtz problems. To alleviate memory and computational complexity, Lucero Lorca et al. (2024) use an iterative method at the highest-level subproblems. In concurrent work to our own, Kump et al. (2025) approach three-dimensional problems with uniform discretizations using a hybrid GPU-CPU approach by combining the composite spectral collocation method with a two-level sparse direct solver. Fortunato et al. (2021) use the ultraspherical spectral method to discretize triangular or quadrilateral mesh elements and compute solutions over polygonal domains by merging DtN operators. Fortunato (2024) develops a variant of the HPS method which merges DtN and ItI operators to solve PDEs on unstructured meshes of smooth two-dimensional surfaces. Beams et al. (2020) develops an implementation of the HPS algorithm targeting parallel shared-memory computer architectures.

There has been significant interest in the GPU acceleration of (low-order) iterative PDE solvers (Georgescu et al., 2013). Other general-purpose packages, such as MFEM and libCEED, implement high-order iterative solvers with GPU acceleration (Abdelfattah et al., 2021; Kolev et al., 2021). Accelerating these algorithms often requires the rapid application of an extremely sparse system matrix. Applying GPU acceleration to direct solvers (Abdelfattah et al., 2022; Ghysels and Synk, 2022; Li and Demmel, 2003) requires different techniques; the literature has mostly focused on sparse direct solvers which do not employ a nested dissection method.

These sparse direct solvers often have much higher peak memory requirements and heterogeneous computation profiles when compared with iterative PDE solvers.

Other solvers have been designed directly for GPU acceleration. Yesypenko and Martinsson

(2024a,b) designed a composite high-order spectral collocation method and associated sparse direct solver with highly heterogeneous computation patterns, which eases GPU acceleration. This method can solve variable-coefficient 2D problems very quickly; our method solves similar problems, but can also handle three-dimensional problems and interface with automatic differentiation. Developing a GPU-compatible implementation of the solver in Yesypenko and Martinsson (2024b), as well as the implementation of our method, has been greatly eased by the advent of high-performance hardware-accelerated linear algebra frameworks popularized by deep learning, such as PyTorch (Ansel et al., 2024) and JAX (Bradbury et al., 2018). These frameworks are highly efficient for batched linear algebra tasks and implement automatic differentiation capabilities. Both are high-level packages that sit on top of the XLA compiler (Leary and Wang, 2017), which compiles and launches optimized kernels that execute on general-purpose GPUs. There has also been work creating automatic differentiation compatible PDE solvers in JAX, tailored for problems such as synchrotron simulation (Diao et al., 2024), computational mechanics (Xue et al., 2023), and ordinary differential equations (Kidger, 2021).

### 3 Introduction to HPS methods

In this section, we provide an overview of the HPS algorithms used in the paper, with a particular eye on their computational structure and the possibilities for GPU acceleration.<sup>1</sup> Full algorithms are available in Sections A to C. We use different variants of this algorithm for merging different types of Poincaré–Steklov operators. A Poincaré–Steklov operator  $T : g \mapsto h$  maps from one type of boundary data to another. Take, for example, a Dirichlet-to-Neumann (DtN) operator, which maps from Dirichlet data  $g$  on the boundary of  $\Omega$  to Neumann data  $h$  on the same boundary:

$$\begin{aligned} g &= u|_{\partial\Omega}, \\ h &= u_n|_{\partial\Omega}, \end{aligned}$$

where  $u$  satisfies Equation (1). Another example commonly used is an impedance-to-impedance (ItI) operator, which maps “incoming” impedance data  $g$  to “outgoing” impedance data  $h$  (Gillman et al., 2015):

$$\begin{aligned} g &= u_n + i\eta u|_{\partial\Omega}, \\ h &= u_n - i\eta u|_{\partial\Omega}, \end{aligned}$$

where  $u$  satisfies Equation (1). These Poincaré–Steklov operators are linear operators, and we work with their discretization  $\mathbf{T}$ . Throughout the algorithm, we also work with  $\mathbf{g}$  and  $\mathbf{h}$ , vectors of incoming and outgoing boundary data evaluated at a set of discretization points.

It is important to remember that because we are solving a linear partial differential equation, we can decompose the solution  $u(x)$  into a particular solution  $v(x)$  and homogeneous solution  $w(x)$  where  $u(x) = v(x) + w(x)$ . The particular solution  $v(x)$  satisfies

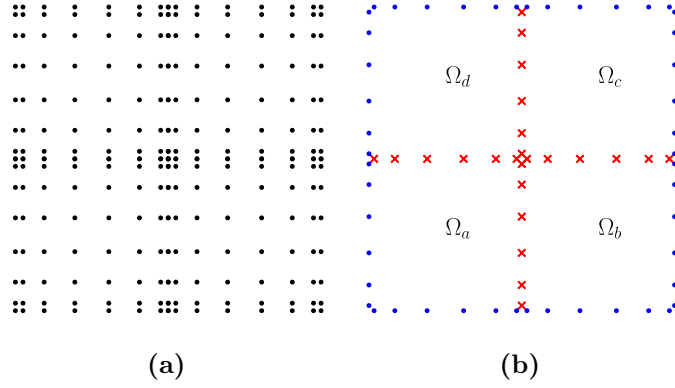
$$\begin{cases} \mathcal{L}v(x) = f(x), & x \in \Omega, \\ v(x) = 0, & x \in \partial\Omega, \end{cases}$$

and the homogeneous solution  $w(x)$  satisfies

$$\begin{cases} \mathcal{L}w(x) = 0, & x \in \Omega, \\ w(x) = g(x), & x \in \partial\Omega. \end{cases}$$

---

<sup>1</sup>For an instructional introduction to these methods, we refer the interested reader to Martinsson (2015); Gillman et al. (2015); Martinsson (2019).



**Figure 1:** Visualizing the high-order composite spectral collocation scheme for a simple two-dimensional problem. Figure 1a shows the Chebyshev points on a two-dimensional problem with polynomial order  $p = 8$  and  $L = 1$  level of refinement. Figure 1b shows the Gauss-Lobatto points discretizing the boundaries of the leaves using order  $q = 6$ . When merging nodes together, it is important to distinguish between the *exterior* boundary points, drawn with blue dots, and the *interior* boundary points, drawn with red  $\times$ 's.

### 3.1 Discretization via composite high-order spectral collocation

To numerically solve Equation (1), a discretization is needed to represent  $\mathcal{L}$ ,  $f$ , and  $g$  in some finite-dimensional basis. HPS methods perform this discretization in two steps: a recursive partition of the domain  $\Omega$  and a high-order spectral collocation scheme. The first step is to recursively partition  $\Omega$  using a quadtree or octree structure down to a user-specified maximum depth  $L$ .

We use  $n_{\text{leaves}}$  to denote the number of patches at the finest level of the spatial partition. In this section, we consider uniform discretization trees, so each tree with depth  $L$  will have  $n_{\text{leaves}} = 2^{dL}$  elements at the lowest level in dimension  $d = 2, 3$ . In Section 4.2, we consider more general discretization trees. At times, it will be useful to describe the progress of the algorithm using language to describe the trees representing the spatial partition. To that end, we will sometimes refer to the elements as *nodes* and elements at the lowest level of the tree as *leaves*. The element at the highest level of the tree, which represents the entire computational domain, is sometimes called the *root*.

Each leaf is discretized using a tensor product of Chebyshev-Lobatto points, with user-specified order  $p$ . This requires  $p^d$  points per leaf in dimension  $d = 2, 3$ . HPS methods typically call for an order- $q$  Gauss-Legendre quadrature rule to represent the boundary of each leaf. For simplicity and stability, we always use  $q = p - 2$ , following Gillman et al. (2015). In two and three dimensions, there are  $(2d)q^{d-1}$  boundary discretization points per leaf. We show the interior and boundary points in Figure 1. The resulting discretization has  $N = n_{\text{leaves}}p^d = (2^L p)^d$  interior discretization points.

### 3.2 Local solve stage

The first step is to discretize the differential operator restricted to that leaf on the  $p^d$  Chebyshev discretization points. We call the resulting matrix  $\mathbf{L}^{(i)}$  for leaf  $i$ . This matrix is a combination of Chebyshev spectral differentiation matrices (Trefethen, 2000) and evaluations of the spatially varying coefficient functions. The source function  $f$  is discretized on the same points, and we call the resulting vector  $\mathbf{f}^{(i)}$ . At this point, the HPS algorithm solves a local boundary-value problem on each patch. Standard practice (Gillman et al., 2015; Fortunato, 2024) is to

solve this problem using a “boundary bordering” technique which enforces the differential operator on the Chebyshev nodes interior to the leaf, and enforces a Dirichlet or impedance boundary condition on the Chebyshev nodes on the leaf’s boundary. At this point in the algorithm, the correct boundary values to enforce at each leaf are unknown, so  $\mathbf{L}^{(i)}$  and  $\mathbf{f}^{(i)}$  are used to precompute a solution map for the local problem. To precompute this local solution map, the algorithm constructs a matrix  $\mathbf{Y}^{(i)}$  which maps from any boundary data  $\mathbf{g}^{(i)}$  to the corresponding homogeneous solution on the Chebyshev nodes. The algorithm also computes  $\mathbf{v}^{(i)}$ , a vector evaluating the particular solution on the Chebyshev nodes. Both  $\mathbf{Y}^{(i)}$  and  $\mathbf{v}^{(i)}$  can be expressed simply in terms of the input data; lines 3 and 4 in Algorithm 7 in the Appendix give these expressions. For each leaf, the algorithm also constructs a Poincaré–Steklov matrix  $\mathbf{T}^{(i)}$  and a vector of outgoing data  $\mathbf{h}^{(i)}$ . Computing  $\mathbf{T}^{(i)}$  and  $\mathbf{h}^{(i)}$  only requires multiplying  $\mathbf{Y}^{(i)}$  and  $\mathbf{v}^{(i)}$  with a fixed, precomputed operator which composes interpolation matrices from Chebyshev to Gauss–Legendre discretization points and spectral Chebyshev differentiation matrices.

Algorithm 1 shows that this stage of the algorithm is a long loop over linear algebra operations. The size of these operations is controlled by the polynomial order  $p$ , and we find these operations are efficient for the orders  $p \leq 16$  considered in this work. The units of work inside the loop are *embarrassingly parallel*, meaning that one iteration does not depend on the output of any other iterations. Furthermore, because we hold  $p$  constant on all leaves, all of the linear algebra operations are homogeneous, meaning they all operate on the same sizes of matrices. This computational structure facilitates GPU acceleration by batching and parallelizing local solves. We give full details describing the local solve stage in Algorithms 7 and 8 in the Appendix.

---

**Algorithm 1:** Local solve stage. Full details are available in Algorithms 7 and 8.

---

**Input:** Discretized differential operators  $\{\mathbf{L}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$ ; discretized source functions  $\{\mathbf{f}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$ ; precomputed interpolation and differentiation matrices

- 1 **for** Leaf  $i = 1, \dots, n_{\text{leaves}}$  **do**
- 2   Perform boundary bordering to  $\mathbf{L}^{(i)}$
- 3   Invert the resulting matrix // Main computational work
- 4   Construct  $\mathbf{Y}^{(i)}$ , the interior solution matrix
- 5   Construct  $\mathbf{T}^{(i)}$ , the Poincaré–Steklov matrix
- 6   Construct  $\mathbf{v}^{(i)}$ , the leaf-level particular solution
- 7   Construct  $\mathbf{h}^{(i)}$ , the outgoing boundary data

**Result:** Poincaré–Steklov matrices  $\{\mathbf{T}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$ ; outgoing boundary data  $\{\mathbf{h}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$ ; interior solution matrices  $\{\mathbf{Y}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$ ; leaf-level particular solutions  $\{\mathbf{v}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$

---

### 3.3 Merge stage

After computing  $\mathbf{T}^{(i)}$  and  $\mathbf{h}^{(i)}$  for each leaf in the local solve stage, the HPS algorithm begins merging nodes of the tree together. This process creates a hierarchy of solution operators which will later be used to propagate the boundary data on  $\partial\Omega$  to the boundary of each leaf. For 2D problems, our implementation merges nodes four at a time, and for 3D problems, our implementation merges nodes eight at a time. Suppose the algorithm is merging a set of nodes  $\{a, b, \dots\}$  which all share parent node  $j$ . The algorithm has access to the following data:

- $\{\mathbf{T}^{(a)}, \mathbf{T}^{(b)}, \dots\}$ , the Poincaré–Steklov matrices of the nodes being merged.



- $\{\mathbf{h}^{(a)}, \mathbf{h}^{(b)}, \dots\}$ , the outgoing boundary data due to the particular solution of the nodes being merged.

At this point, it is helpful to distinguish between vectors that are defined along the *exterior* of the patches being merged and vectors that are defined along the *interior* of the patches being merged. We indicate these vectors with subscripts *ext* and *int*, respectively. See Figure 1b for a diagram of the interior and exterior points in a 2D merge operation. The goal of the merge operation is to precompute a solution operator which propagates the information from the exterior boundary points to the interior boundary points. This solution operator takes the form  $\mathbf{g}_{\text{ext}}^{(j)} \mapsto \mathbf{S}^{(j)} \mathbf{g}_{\text{ext}}^{(j)} + \tilde{\mathbf{g}}^{(j)}$ . In this equation,  $\mathbf{S}^{(j)} \mathbf{g}_{\text{ext}}^{(j)}$  evaluates the homogeneous solution on the interior boundary points, and  $\tilde{\mathbf{g}}^{(j)}$  evaluates the particular solution on the interior boundary points. As in Section 3.2, the boundary data  $\mathbf{g}_{\text{ext}}^{(j)}$  is not available at this stage of the algorithm, but it is possible to precompute the other parts of the solution operator. To that end, each merge operation will compute:

- $\mathbf{S}^{(j)}$ , the propagation operator for node  $j$ , which maps incoming homogeneous boundary data from the exterior boundary points to the interior boundary points.
- $\tilde{\mathbf{g}}^{(j)}$ , the incoming boundary data due to the particular solution evaluated at the interior boundary points.
- $\mathbf{T}^{(j)}$ , the Poincaré–Steklov matrix for node  $j$ .
- $\mathbf{h}^{(j)}$ , the outgoing boundary data for node  $j$ .

$\mathbf{S}^{(j)}$  and  $\tilde{\mathbf{g}}^{(j)}$  will be used in the final stage of the HPS method when applying the precomputed solution operator, and  $\mathbf{T}^{(j)}$  and  $\mathbf{h}^{(j)}$  will be used in a future merge operation.

To compute the merge outputs, the HPS algorithm sets up a system of equations for a given  $\mathbf{g}_{\text{ext}}^{(j)}$  and unknown  $\mathbf{g}_{\text{int}}^{(j)}$  and solve using a Schur complement approach.<sup>2</sup> The constraints in this system come from our knowledge that the solution and its derivative will be continuous across merge interfaces. The resulting system of constraints is a linear system and can be written in a blockwise fashion:

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} \begin{bmatrix} \mathbf{g}_{\text{ext}}^{(j)} \\ \mathbf{g}_{\text{int}}^{(j)} \end{bmatrix} = \begin{bmatrix} \mathbf{u}_{\text{ext}}^{(j)} - \mathbf{h}_{\text{ext}}^{(\text{child})} \\ -\mathbf{h}_{\text{int}}^{(\text{child})} \end{bmatrix}. \quad (3)$$

The next step is to build the matrices  $\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}$  blockwise from the child Poincaré–Steklov matrices  $\{\mathbf{T}^{(a)}, \mathbf{T}^{(b)}, \dots\}$ , build  $\mathbf{h}_{\text{ext}}^{(\text{child})}$  and  $\mathbf{h}_{\text{int}}^{(\text{child})}$  from the child outgoing boundary data  $\{\mathbf{h}^{(a)}, \mathbf{h}^{(b)}, \dots\}$ , and use  $\mathbf{u}_{\text{ext}}^{(j)}$  to represent the unknown solution evaluated on the exterior boundary points. For 2D DtN merges, these objects are defined in Equations (22) to (28). Derivations for the 2D ItI and 3D DtN cases can be found in Sections B and C. At this point in the algorithm,  $\mathbf{u}_{\text{ext}}^{(j)}$  is unknown, which means one can not directly invert the linear system to solve for  $\mathbf{g}_{\text{ext}}^{(j)}$  or  $\mathbf{g}_{\text{int}}^{(j)}$ . However, one can use a Schur complement approach to partially solve the system:

$$\begin{bmatrix} \mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C} & \mathbf{0} \\ \mathbf{D}^{-1}\mathbf{C} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{g}_{\text{ext}}^{(j)} \\ \mathbf{g}_{\text{int}}^{(j)} \end{bmatrix} = \begin{bmatrix} \mathbf{u}_{\text{ext}}^{(j)} - \mathbf{h}_{\text{ext}}^{(\text{child})} - \mathbf{B}\mathbf{D}^{-1}\mathbf{h}_{\text{int}}^{(\text{child})} \\ -\mathbf{D}^{-1}\mathbf{h}_{\text{int}}^{(\text{child})} \end{bmatrix}.$$

---

<sup>2</sup>To the best of our knowledge, the block system associated with the merge four box procedure for the 2D DtN method was first documented in Chipman et al. (2024). We generalize this presentation to encompass 3D problems and merging ItI matrices as well.

Interpreting the rows of this linear system gives us the desired outputs:

$$\mathbf{u}_{\text{ext}}^{(j)} = \underbrace{(\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})}_{\mathbf{T}^{(j)}} \mathbf{g}_{\text{ext}}^{(j)} + \underbrace{\mathbf{h}_{\text{ext}}^{(\text{child})} - \mathbf{B}\mathbf{D}^{-1}\mathbf{h}_{\text{int}}^{(\text{child})}}_{\mathbf{h}^{(j)}}, \quad (4)$$

$$\mathbf{g}_{\text{int}}^{(j)} = \underbrace{-\mathbf{D}^{-1}\mathbf{C}}_{\mathbf{S}^{(j)}} \mathbf{g}_{\text{ext}}^{(j)} + \underbrace{-\mathbf{D}^{-1}\mathbf{h}_{\text{int}}^{(\text{child})}}_{\tilde{\mathbf{g}}^{(j)}}. \quad (5)$$

Algorithm 2 gives pseudocode for the merge stage of the HPS algorithm. The majority of the computational work for each merge is inverting  $\mathbf{D}$ , which has size proportional to the number of discretization points along the merge interfaces. This matrix is quite small at the lowest levels of the discretization tree and grows as the algorithm proceeds to higher nodes in the tree. Similar to the local solve stage, each merge operation is dominated by linear algebra work, and the units of work inside the inner loop are embarrassingly parallel. This means we can easily use GPU acceleration to parallelize the inner loop of Algorithm 2. The outer loop is iterating over different levels, which means the computation during outer loop iteration  $\ell$  depends on the outputs of the previous iteration. Because we only use a moderate number of refinement levels  $L < 10$ , we find Algorithm 2 executes very quickly on the GPU despite this dependency structure. We note that our choice to merge nodes four-to-one and eight-to-one (rather than the standard two-to-one) decreases the length of the outer loop by factors of two and three, respectively.

---

**Algorithm 2:** Merge stage. Full details are available in Sections A.2, B.2 and C.2.

---

**Input:** Leaf-level Poincaré–Steklov matrices  $\{\mathbf{T}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$ ; Leaf-level outgoing boundary data  $\{\mathbf{h}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$

- 1 **for** Merge level  $\ell = L - 1, \dots, 0$  **do**
- 2   **for** Node  $j$  in level  $\ell$  **do**
- 3     Let  $a, b, \dots$  be the children of node  $j$
- 4     Use  $\{\mathbf{T}^{(a)}, \mathbf{T}^{(b)}, \dots\}$  to build blocks  $\mathbf{A}, \mathbf{B}, \mathbf{C}$ , and  $\mathbf{D}$
- 5     Use  $\{\mathbf{h}^{(a)}, \mathbf{h}^{(b)}, \dots\}$  to build  $\mathbf{h}_{\text{ext}}^{(\text{child})}$  and  $\mathbf{h}_{\text{int}}^{(\text{child})}$
- 6     Invert  $\mathbf{D}$  // Main computational work
- 7     Evaluate  $\mathbf{T}^{(j)}, \mathbf{h}^{(j)}, \mathbf{S}^{(j)}, \tilde{\mathbf{g}}^{(j)}$  // Equations (4) and (5)

**Result:** Poincaré–Steklov matrices  $\mathbf{T}^{(j)}$  for each node; outgoing boundary data  $\mathbf{h}^{(j)}$  for each node; propagation operators  $\mathbf{S}^{(j)}$  for each node; incoming particular solution data  $\tilde{\mathbf{g}}^{(j)}$  for each node

---

### 3.4 Downward pass

In the final stage of the HPS method, all parts of the structured solution operators mapping  $g \mapsto u$  have been computed. The HPS algorithm evaluates this structured solution operator by propagating information down the discretization tree from the boundary of the root to the interior of the leaves. The  $\mathbf{S}^{(j)}$  matrices propagate the homogeneous boundary data to the merge interfaces, and the  $\tilde{\mathbf{g}}^{(j)}$  vectors add back in the particular solution (line 4 in Algorithm 3.) This part of the HPS method is extremely fast, as it only involves matrix-vector products. As with the structure of the merge stage, the iterations of the inner loop are embarrassingly parallel and can be batched on the GPU. We show the pseudocode for this stage in Algorithm 3.



---

**Algorithm 3:** Downward pass.

---

**Input:** Boundary data  $\mathbf{g}$ ; propagation operators  $\mathbf{S}^{(j)}$  for each node; incoming particular solution data  $\tilde{\mathbf{g}}^{(j)}$  for each node; leaf-level interior solution matrices  $\{\mathbf{Y}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$ ; leaf-level particular solutions  $\{\mathbf{v}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$

- 1 **for** Merge level  $\ell = 0, \dots, L - 1$  **do**
- 2     **for** Node  $j$  in level  $\ell$  **do**
- 3         Look up  $\mathbf{S}^{(j)}$ ,  $\mathbf{g}^{(j)}$ , and  $\tilde{\mathbf{g}}^{(j)}$
- 4          $\mathbf{g}_{\text{int}} = \mathbf{S}^{(j)}\mathbf{g}^{(j)} + \tilde{\mathbf{g}}^{(j)}$
- 5         Let  $a, b, \dots$  be the children of node  $j$
- 6         Concatenate  $\mathbf{g}_{\text{int}}$  and  $\mathbf{g}^{(j)}$  to form  $\{\mathbf{g}^{(a)}, \mathbf{g}^{(b)}, \dots\}$
- 7 **for** Leaf  $i = 1, \dots, n_{\text{leaves}}$  **do**
- 8      $\mathbf{u}^{(i)} = \mathbf{Y}^{(i)}\mathbf{g}^{(i)} + \mathbf{v}^{(i)}$

**Result:** Leaf-level solutions on the Chebyshev discretization points  $\{\mathbf{u}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$

---

## 4 Hardware acceleration for HPS methods

While the algorithms presented in Section 3 have attractive computational complexity ( $O(p^6 n_{\text{leaves}} + p^3 n_{\text{leaves}}^{3/2})$  in 2D and  $O(p^9 n_{\text{leaves}} + p^6 n_{\text{leaves}}^2)$  in 3D) and possibilities for parallel execution, they also incur large memory footprints. At each step of the algorithm, dense solution matrices are precomputed and must be stored for future use. The outputs of the merge stage dominate the memory complexity of the method. In 2D, the overall memory complexity of storing these matrices is  $O(p^2 n_{\text{leaves}} L)$ , and in 3D, the memory complexity is  $O(p^4 n_{\text{leaves}} 2^L)$ , with a large prefactor. This poses a significant challenge for GPU acceleration as general-purpose GPU architectures have significantly more processor cores per unit of memory than standard multicore CPU nodes. While standard multicore compute nodes may have 1TB of available random-access memory (host RAM), high-end GPUs have only 80GB of on-device memory, with slow interconnects between the GPU and host RAM. This means that batched linear algebra operations are extremely fast on the GPU, but the overall algorithm is slowed by steps transferring data between the GPU and the host. Thus, to efficiently accelerate HPS algorithms on the GPU, one must devote significant thought to reducing the memory footprint of these algorithms. In this section, we introduce two ideas to reduce this memory footprint in two and three-dimensional problems.

### 4.1 Recomputation strategies to minimize communication costs

CPU-bound implementations of HPS methods in 2D often spend an order of magnitude longer in the local solve stage than in the merge or downward pass stages (Fortunato, 2024). This suggests that HPS schemes can be greatly accelerated by placing the local solve stage computation on the GPU alone. Indeed, such savings have been observed in Yesypenko and Martinsson (2024a). We observe that for the lowest levels of the merge stage, each unit of computational work is similarly small, suggesting that the GPU can efficiently accelerate this part of the algorithm, too, provided the algorithm is correctly expressed to leverage its inherent parallelism. In many GPUs, the interconnect between the host device and the GPU’s on-device memory is very slow in relation to the speed at which the thousands of processor cores can process data. This means that for large problem sizes, operations transferring precomputed solution operators to and from the GPU are a major impediment to fast execution as they incur a large latency and are often “blocking” operations, which require all parallel threads to complete before executing.

For large problem sizes in 2D, it is advantageous to delete some data computed in the early stages of the algorithm and recompute it later when necessary. This strategy increases the number of floating point operations on the GPU but minimizes costly data transfers. A leaf-level recomputation strategy for implementing HPS algorithms on a GPU is presented in Algorithm 4; a similar method is presented in Yesypenko and Martinsson (2024a). The leaf recomputation strategy avoids transferring the  $\{\mathbf{Y}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$  and  $\{\mathbf{v}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$  by performing the local solve stage again at the end of the algorithm. Under this recomputation strategy, all of the leaf-level Poincaré–Steklov matrices must be transferred to RAM during the local solve stage, and then back to the GPU during the merge stage.

We find that it is advantageous to push the idea of reducing data transfers at the cost of more floating-point operations further. In our proposed recomputation strategy (Algorithm 5), we delete and recompute the products of the local solve stage and multiple levels of the merge stage. To implement this, we operate in batches defined by “complete subtrees”, which are subtrees containing all of the descendants of a particular node  $j$ . We break the lowest levels of the discretization tree into the largest complete subtrees where the computations in Algorithms 1 and 2 can all fit into a GPU’s on-device memory. The size of these maximal complete subtrees varies depending on GPU memory, polynomial order  $p$ , and floating-point datatype; in our experiments, these maximal complete subtrees usually have depth 6 or 7. For each such complete subtree, we perform all of the local solve and merge operations, after which we only save the top-level Poincaré–Steklov matrix and outgoing boundary data vector. Because this is a small amount of data, we can store it on the GPU, and do not need to move the outputs to host RAM. After processing all of the subtrees, the final merge stages are performed on the GPU. The downward pass is evaluated sequentially on the different subtrees, at which point the local solve stage and low-level merges must be recomputed. This recomputation method was inspired by optimizations for contemporary deep learning architectures (Dao et al., 2022; Gu and Dao, 2024), which suggest kernel fusion, a technique that performs multiple steps of a sequential computation at once to keep the necessary data near the processor cores. We visualize the different recomputation methods in Figure 2.

In Figure 3, we compare the performance of our method across computer architectures and recomputation strategies. We consider two different architectures, a multicore Intel Xeon node with a 64-core processor, and a GPU architecture using a single Nvidia H100 GPU. Evaluating our method on the GPU gives us significant speedups over the multicore CPU architecture, even when using an implementation with no recomputation strategy, which transfers all precomputed matrices to host RAM after each algorithm step. The two recomputation strategies begin to diverge for problem sizes over  $10^7$  discretization points, at which point the precomputed matrices cannot all fit on the GPU. We use subtree depth 7 for this experiment; we explore the effect of this choice in Section E by measuring the runtime for different subtree depths. We also estimate the percentage of peak double-precision floating point operations per second (FLOPS) achieved by the different recomputation strategies. Our proposed recomputation strategy uses the most floating-point operations and has the fastest runtime, which means it reaches a higher percentage of peak FLOPS than the other implementations.

## 4.2 An adaptive discretization strategy to reduce memory complexity in 3D

When extending from two to three dimensions, we face different computational challenges. A large part of the difficulty involves the size of the matrices arising in the merge stage discussed in Section 3.3. For each merge operation, the matrix  $\mathbf{D}$  must be inverted. This matrix has a number of rows and columns proportional to the number of discretization points that lie along the interfaces being merged. In two dimensions, the size of this merge interface is  $O(p2^\ell)$  to

---

**Algorithm 4:** Leaf recomputation strategy.

---

**Input:** Differential operators  $\{\mathbf{L}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$ ; source functions  $\{\mathbf{f}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$ ; boundary data  $\mathbf{g}$

- 1 Let  $b$  be the maximum batch size that can fit on the GPU
- 2 Split the indices  $\{1, 2, \dots, n_{\text{leaves}}\}$  into batches  $I_1, I_2, \dots, I_{\lceil n_{\text{leaves}}/b \rceil}$
- 3 **for** Batch  $j$  **do**
- 4   Move  $\{\mathbf{L}^{(i)}\}_{i \in I_j}$  and  $\{\mathbf{f}^{(i)}\}_{i \in I_j}$  to the GPU
- 5   Perform the local solve stage for this batch of leaves
- 6   Delete  $\{\mathbf{Y}^{(i)}\}_{i \in I_j}$  and  $\{\mathbf{v}^{(i)}\}_{i \in I_j}$
- 7   Transfer  $\{\mathbf{T}^{(i)}\}_{i \in I_j}$  and  $\{\mathbf{h}^{(i)}\}_{i \in I_j}$  to host RAM
- 8 Concatenate  $\{\mathbf{T}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$  and  $\{\mathbf{h}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$  and transfer to GPU
- 9 Perform all merge operations on the GPU and transfer all  $\mathbf{S}^{(i)}$  and  $\tilde{\mathbf{g}}^{(i)}$  to host RAM
- 10 Propagate boundary data to the leaves
- 11 Transfer leaf-level boundary data  $\{\mathbf{g}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$  to host RAM
- 12 **for** Batch  $j$  **do**
- 13   Move  $\{\mathbf{L}^{(i)}\}_{i \in I_j}$ ,  $\{\mathbf{f}^{(i)}\}_{i \in I_j}$ , and  $\{\mathbf{g}^{(i)}\}_{i \in I_j}$  to the GPU
- 14   Compute local solutions  $\mathbf{u}_{i \in I_j}^{(i)}$
- 15   Transfer  $\{\mathbf{u}^{(i)}\}_{i \in I_j}$  to host RAM

**Result:** Solutions on the Chebyshev discretization points  $\{\mathbf{u}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$

---

---

**Algorithm 5:** Subtree recomputation strategy.

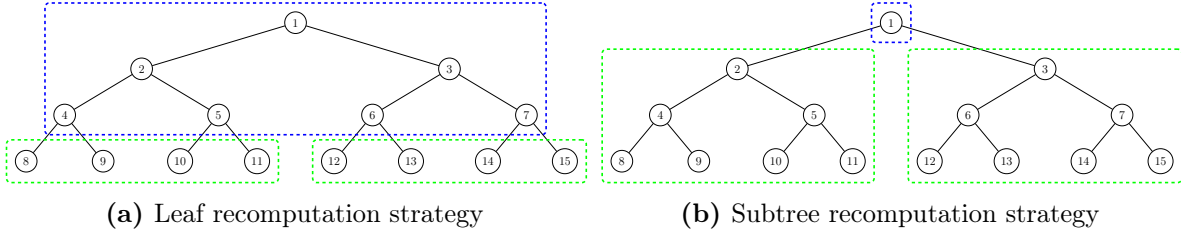
---

**Input:** Differential operators  $\{\mathbf{L}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$ ; source functions  $\{\mathbf{f}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$ ; boundary data  $\mathbf{g}$

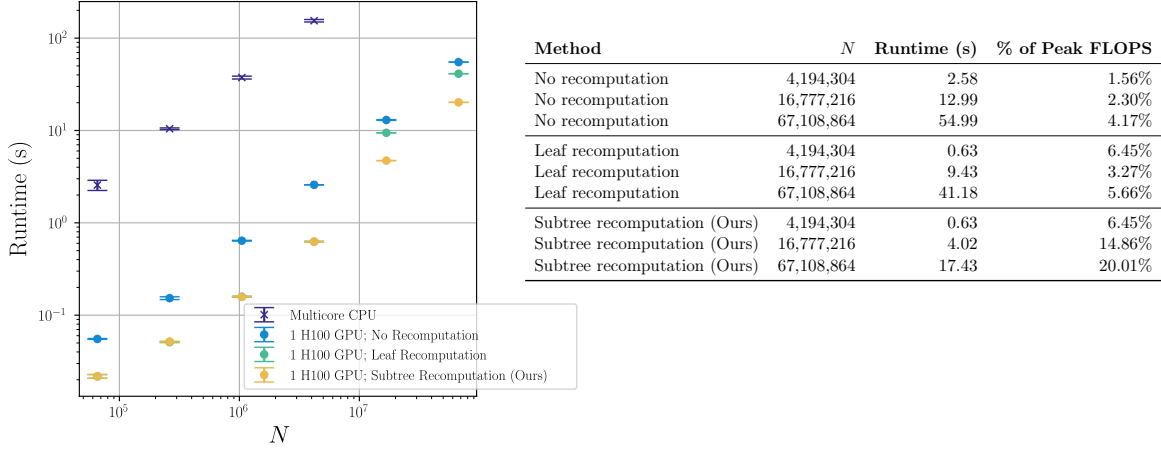
- 1 Let  $M = m_1, m_2, \dots$  be the roots of the maximal subtrees
- 2 **for** Subtree rooted at  $m_j$  **do**
- 3   Let  $I_j$  be the set of leaves of the subtree
- 4   Move  $\{\mathbf{L}^{(i)}\}_{i \in I_j}$  and  $\{\mathbf{f}^{(i)}\}_{i \in I_j}$  to the GPU
- 5   Perform the local solve stage for this subtree
- 6   Delete  $\{\mathbf{Y}^{(i)}\}_{i \in I_j}$  and  $\{\mathbf{v}^{(i)}\}_{i \in I_j}$
- 7   Merge the leaves to the top of subtree  $m_j$ , deleting all outputs except  $\mathbf{T}^{(m_j)}$  and  $\mathbf{h}^{(m_j)}$
- 8   Keep  $\mathbf{T}^{(m_j)}$  and  $\mathbf{h}^{(m_j)}$  on GPU
- 9 Perform final merge operations on the GPU
- 10 Propagate boundary data to the roots of the maximal subtrees
- 11 Transfer boundary data to host RAM
- 12 **for** Subtree rooted at  $m_j$  **do**
- 13   Let  $I_j$  be the set of leaves of the subtree
- 14   Move  $\{\mathbf{L}^{(i)}\}_{i \in I_j}$ ,  $\{\mathbf{f}^{(i)}\}_{i \in I_j}$  and  $\mathbf{g}^{(m_j)}$  to the GPU
- 15   Perform the local solve stage for this subtree
- 16   Merge the leaves to the top of subtree  $j$
- 17   Propagate boundary information down the tree to the leaves
- 18    $\mathbf{u}^{(i)} \leftarrow \mathbf{Y}^{(i)}\mathbf{g}^{(i)} + \mathbf{v}^{(i)}$
- 19   Transfer  $\{\mathbf{u}^{(i)}\}_{i \in I_j}$  to host RAM

**Result:** Solutions on the Chebyshev discretization points  $\{\mathbf{u}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$

---



**Figure 2:** Comparing the batching patterns of the two different recomputation strategies. The leaf recomputation strategy in Figure 2a performs the local solve stage operations in large batches and then performs all of the merge stages in a separate batch. Our proposed subtree recomputation strategy (Figure 2b) performs local solves and multiple levels of the merge stage for a complete subtree of the discretization tree structure.



**Figure 3:** Even with a naïve implementation of the HPS algorithm which does not perform any recomputation, using a single GPU achieves large speedups over a multicore CPU system. When we use our proposed subtree recomputation strategy, the speedup increases by another factor of two. (Left) We vary  $L = 4, \dots, 9$  and hold  $p = 16$  fixed to generate problems with  $N = p^2 4^L = 256 \times 4^L = 4^{L+4}$  degrees of freedom. We measure the total runtime of our 2D method merging DtN matrices; this runtime includes the execution of the local solve stage, the merge stage, and the downward pass. Vertical error bars show  $\pm 1$  standard error computed over five trials. (Right) For each of the GPU implementations, we compute the total number of FLOPS and report this as a percentage of the GPU’s peak FLOPS, estimated by the manufacturer to be  $34 \times 10^{12}$ .

merge nodes  $\ell$  levels above the leaves. In three dimensions, the size of this merge interface is  $O(p^2 4^\ell)$ . Because this quantity grows very quickly for 3D problems as we increase the tree depth  $L$ , we quickly run out of memory required to store and invert the matrices on the GPU at the highest level of the merge stage. Performing the top-level merge operation is when the instantaneous memory footprint peaks, as space for  $\mathbf{D}$ ,  $\mathbf{D}^{-1}$ , and various buffers must all be allocated on the GPU simultaneously. In contrast with other parts of the algorithm, this peak memory footprint is not reducible by strategies such as batching or data transfer. Because of this need to invert matrices near the memory limits of the GPU, we transfer data to and from the GPU at each merge level and do not use the recomputation strategies discussed in Section 4.1.

To reduce the size of  $\mathbf{D}$  at the final merge step, we propose to extend the adaptive HPS method presented for 2D problems in Geldermans and Gillman (2019) to three dimensions.

This method adaptively refines element sizes in a data-dependent manner, in an effort to concentrate discretization points in the regions of the domain where the coefficient and source functions have high local variation. In Table 1, we show that the size of  $\mathbf{D}$  generated using our adaptive refinement technique is much smaller than that of the uniform refinement with no loss in accuracy.

Developing a version of the HPS methods presented in Section 3 that is compatible with an adaptive discretization requires slight modification of the algorithms presented in the previous section. The major changes are the introduction of a method for adaptively refining our octree and a method for merging nodes with different levels of refinement.

Method	$p$	Relative $\ell_\infty$ Error	# Leaves	Size of $\mathbf{D}$
Uniform	8	$1.48 \times 10^{-4}$	512	6,912
Adaptive	8	$1.45 \times 10^{-4}$	190	2,700
Uniform	12	$3.62 \times 10^{-7}$	512	19,200
Adaptive	12	$2.04 \times 10^{-7}$	442	7,500
Uniform	16	$4.20 \times 10^{-6}$	64	9,408
Adaptive	16	$1.41 \times 10^{-6}$	57	4,116

**Table 1:** Adaptive discretization methods can greatly reduce the peak memory requirements of HPS methods in three dimensions. We present the size of the discretization tree and final merge steps in our 3D “wavefront” example (Section 6.1). We compare adaptive and uniform discretizations that have similar errors and observe the adaptive discretization strategy can greatly reduce the size of the final  $\mathbf{D}$  matrix, a proxy for peak memory usage.

#### 4.2.1 Criterion for adaptive refinement

For a given leaf of the discretization tree, let  $\mathbf{x}_0$  be the set of  $p^3$  Chebyshev points discretizing the leaf. Let  $\mathbf{x}_1$  be the set of  $8p^3$  discretization points found by breaking the leaf into eight children and creating a Chebyshev grid on each child. Let  $\mathbf{L}_{8f1}$  be an interpolation matrix mapping from  $\mathbf{x}_0$  to  $\mathbf{x}_1$ . We evaluate whether a function is sufficiently refined on a leaf by checking whether we can use polynomial interpolation to accurately map from evaluations on  $\mathbf{x}_0$  to evaluations on  $\mathbf{x}_1$ , relative to the global  $L_\infty$  norm of the function. We specify a tolerance parameter  $\epsilon$ , and for each leaf in our tree, we check the following condition:

$$\frac{\|f(\mathbf{x}_1) - \mathbf{L}_{8f1}f(\mathbf{x}_0)\|_\infty}{\|f\|_\infty} < \epsilon. \quad (6)$$

If this condition is met, we say the leaf is sufficiently refined. Otherwise, we split the leaf into eight children and check each child. We form a final discretization tree by refining each coefficient function in our differential operator, as well as the source term, and taking the union of the resulting trees. Additionally, we refine a few extra leaves to enforce a “level restriction” criterion, which specifies that no leaf can have a side length greater than twice that of its neighbors. This method is an extension of the method for two-dimensional problems presented in Geldermans and Gillman (2019), which uses a similar relative  $L_2$  convergence criterion and level restriction criterion.

#### 4.2.2 Local solve stage

The local solve stage for adaptively refined discretization trees is the same as in the uniform refinement case. Although they are defined over leaves with different volumes, each local

boundary value problem has the same number of interior and boundary discretization points, and the local problems are still embarrassingly parallel. Thus, we can use batched linear algebra to accelerate this part of the algorithm.

#### 4.2.3 Merging nodes with different discretization levels

The nonuniform merge stage is different from the uniform merge stage because neighboring nodes may have different refinement levels, which means the discretization points along either side of the merge interface may not exactly align. Recall the block linear system (Equation (3)) arising during the merge stage:

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} \begin{bmatrix} \mathbf{g}_{\text{ext}}^{(j)} \\ \mathbf{g}_{\text{int}}^{(j)} \end{bmatrix} = \begin{bmatrix} \mathbf{u}_{\text{ext}}^{(j)} - \mathbf{h}_{\text{ext}}^{(\text{child})} \\ -\mathbf{h}_{\text{int}}^{(\text{child})} \end{bmatrix}.$$

When neighboring volume elements have different refinement levels, there will be a mismatch between the interior boundary discretization points on either side of the merge interface. We need to decide how to represent  $\mathbf{g}_{\text{int}}^{(j)}$ ,  $\mathbf{h}_{\text{int}}^{(\text{child})}$ ,  $\mathbf{B}$ ,  $\mathbf{C}$ , and  $\mathbf{D}$  in Equation (3). We choose to discretize these objects using the coarser of the two sets of discretization points along the merge interface; the discretization points along the exterior boundary elements are inherited from the child nodes. To assemble the blocks in Equation (3), this requires projecting some rows and columns of the Poincaré–Steklov matrices using precomputed interpolation operators which map between one and four 2D Gauss–Legendre panels. The “level restriction” constraint greatly simplifies this compression because the resulting projection operations are guaranteed to be four-to-one.

#### 4.2.4 Downward pass

To propagate the boundary information to the leaf nodes, we follow the general structure of Algorithm 3. However, we must undo the projection along merge interfaces that occurs during the nonuniform merge stage. This is accomplished by applying the precomputed interpolation operators to the boundary data  $\mathbf{g}^{(i)}$ .

## 5 Numerical examples in two dimensions

In this section, we present numerical results on problems with two spatial dimensions. All experiments in this section were conducted using one Nvidia H100 GPU and a host memory space with 100GB of RAM. In all of the experiments, we use the novel subtree recomputation strategy introduced in Section 4.1; the DtN version of this recomputation strategy uses subtrees of depth 7, and the ItI version of this recomputation strategy uses subtrees of depth 6.

### 5.1 High-order convergence on variable-coefficient problems with known solutions

There are two main ways to increase the accuracy of our composite spectral collocation scheme: refine each leaf patch into four children or increase the polynomial order of the representation of the solution on each patch. Empirically, the error is controlled by the polynomial order  $p$  and the side length of each leaf  $h$ . In our implementation,  $p$  is specified by the user and  $h$  is controlled by  $L$ , the user-specified depth of the discretization tree. The tradeoff between these two parameters is a widely-studied topic in numerical analysis and goes by the name of “ $hp$ -adaptivity”. We study the  $hp$ -adaptivity properties of our solver using two problems with



variable-coefficient differential operators and known solutions. The first problem is a variant of Poisson's equation with spatially-varying coefficients:

$$\begin{cases} \Delta u(x) - \cos(5x_2)u_{x_1}(x) + \sin(5x_2)u_{x_2}(x) = f(x), & x \in [-1, 1]^2, \\ u(x) = g(x), & x \in \partial[-1, 1]^2, \end{cases} \quad (7)$$

where  $u_{x_1}$  and  $u_{x_2}$  are the partial derivatives of  $u$  in directions  $x_1$  and  $x_2$ , respectively. We manufacture the source  $f$  and the Dirichlet data  $g$  so the solution to this problem is

$$u(x) = u(x_1, x_2) = e^{5x_1} \sin(5x_2) + \sin(10\pi x_1) \sin(\pi x_2).$$

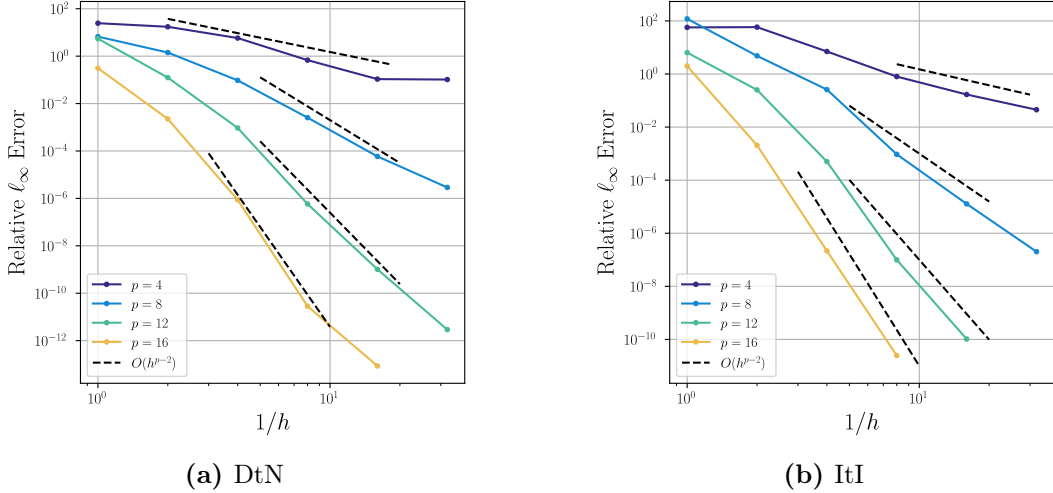
We also study an inhomogeneous Helmholtz problem with a Robin boundary condition:

$$\begin{cases} \Delta u(x) + (1 + e^{-50\|x\|^2})u(x) = f(x), & x \in [-1, 1]^2, \\ u_n(x) + iu(x) = g(x), & x \in \partial[-1, 1]^2. \end{cases} \quad (8)$$

We manufacture the source  $f$  and the Robin data  $g$  so solution to this problem is

$$u(x) = u(x_1, x_2) = e^{i20x_1} + e^{i30x_2}.$$

Figure 4 shows the convergence of our method on these problems. We measure the relative error of our computed solution  $\mathbf{u}$  by computing  $\|\mathbf{u} - \mathbf{u}_{\text{true}}\|_{\infty} / \|\mathbf{u}_{\text{true}}\|_{\infty}$ . The  $\ell_{\infty}$  norms are estimated by taking the maximum over all interior discretization points. In Figure 4, we see the errors of both the DtN and ItI versions of the method converging at rate  $O(h^{p-2})$ , even for high polynomial orders.



**Figure 4:** Using  $p$  Chebyshev points per dimension on each leaf, and leaves of size  $h$ , the relative  $\ell_{\infty}$  errors of our method converge at rate  $O(h^{p-2})$ . Figure 4a shows the convergence of the HPS method using DtN matrices applied to Equation (7) and Figure 4b shows the convergence of the HPS method using ItI matrices applied to Equation (8).

## 5.2 High-frequency forward wave scattering problems

In this example, we solve a variable-coefficient Helmholtz equation coupled with a Sommerfeld radiation condition. The system is excited by a plane wave with direction  $\hat{s} = [1, 0]^T$  and frequency  $k$ :

$$\begin{cases} \Delta u(x) + k^2(1 + q(x))u(x) = -k^2 q(x)e^{ik\langle \hat{s}, x \rangle}, & x \in [-1, 1]^2, \\ \sqrt{r} \left( \frac{\partial u}{\partial r} - iku \right) \rightarrow 0, & r = \|x\|_2 \rightarrow \infty. \end{cases} \quad (9)$$

Equation (9) models time-harmonic wave scattering in many imaging modalities, such as sonar or radar imaging, geophysical sensing, and nondestructive testing of materials (Borges et al., 2017). As such, forward wave scattering solvers are often used inside an inner loop of optimization routines for solving these inverse problems. These forward solvers must be highly optimized as they are evaluated hundreds or thousands of times over the course of an algorithm.

To solve Equation (9), we use the ItI variant of our implementation, and at the top level of the merge stage, we solve a boundary integral equation which enforces the Sommerfeld radiation condition (Gillman et al., 2015). Discretizing the boundary integral equation requires a high-order Nyström method to generate single- and double-layer potential matrices, which we perform in MATLAB using the chunkIE package (Askham et al., 2024). For each discretization level and frequency, generating these matrices takes a few seconds on a standard laptop. Because these matrices can be precomputed once for a given discretization level and frequency, we do not include the time required to generate these matrices in our runtime measurements. The solution of this boundary integral equation specifies incoming impedance data, which is propagated down the tree to form the interior solution. While we solve Equation (9) for one source direction  $\hat{s}$ , this scheme can compute solutions for multiple different sources in parallel at the cost of a few extra matrix-vector multiplications.

In Figures 6 and 8, we measure the runtime and accuracy of our GPU-accelerated solver. Because analytical solutions for Equation (9) are unavailable for general scattering potentials  $q(x)$ , we measure error relative to an overrefined reference solution  $\mathbf{u}_{\text{over}}$  with approximately 2,800 discretization points in both dimensions. For each computed solution  $\mathbf{u}$ , we compute the relative  $\ell_\infty$  error  $\|\mathbf{u} - \mathbf{u}_{\text{over}}\|_\infty / \|\mathbf{u}_{\text{over}}\|_\infty$ . The  $\ell_\infty$  norm is estimated by taking the maximum over a grid of  $500 \times 500$  regularly-spaced grid points. We repeat this experiment for two different choices of the scattering potential  $q(x)$ . We first choose a single Gaussian bump, which loosely focuses the incoming wave:

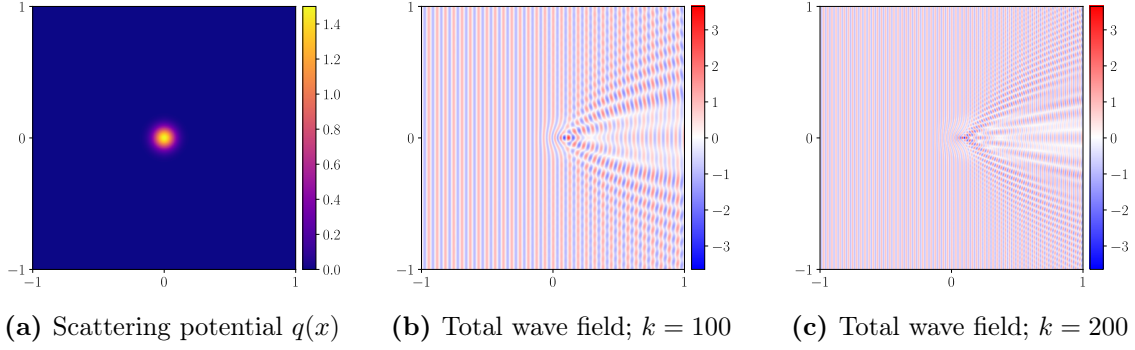
$$q(x) = 1.5e^{-160\|x\|^2}.$$

We also consider a collection of randomly placed Gaussian bumps with centers  $\{z^{(i)}\}_{i=1}^{10}$ , which causes multiple scattering effects at high wavenumbers:

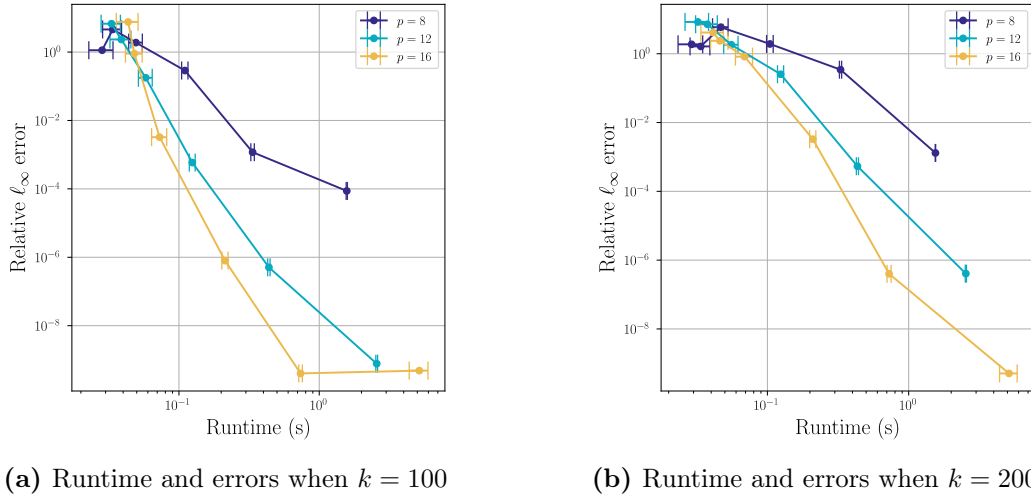
$$q(x) = \sum_{i=1}^{10} e^{-50\|x - z^{(i)}\|^2}. \quad (10)$$

In Figures 5 and 7, we show these scattering potentials as well as the resulting total wave field  $u(x) + e^{ik\langle \hat{s}, x \rangle}$  for different choices of  $k$ .

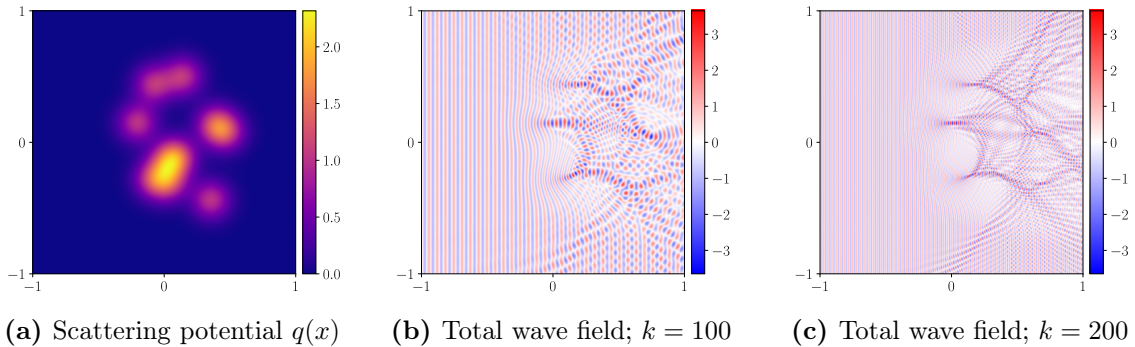
Our GPU-accelerated method with our proposed recomputation strategy is very fast. It is able to compute high-accuracy solutions to these challenging scattering problems in a few seconds. We note that because this method does not rely on any iterative algorithms, the runtime for a fixed discretization level does not depend on the frequency  $k$  or the scattering potential  $q(x)$ . However, finer discretizations are required to achieve a fixed error tolerance as  $k$  increases. Figures 6 and 8 show that polynomial order  $p = 16$  dominates the lower-order methods for all values of  $k$  and scattering potentials considered. In these plots, we vary the number of degrees of freedom  $N = p^2 4^L$  by varying both  $p$  and  $L$ . Different colors correspond to different polynomial orders  $p \in \{8, 12, 16\}$ , and for a given polynomial order, we vary  $L \in \{2, 3, \dots, 7\}$ . This generates problems with degrees of freedom varying between 512 and 4,194,304.



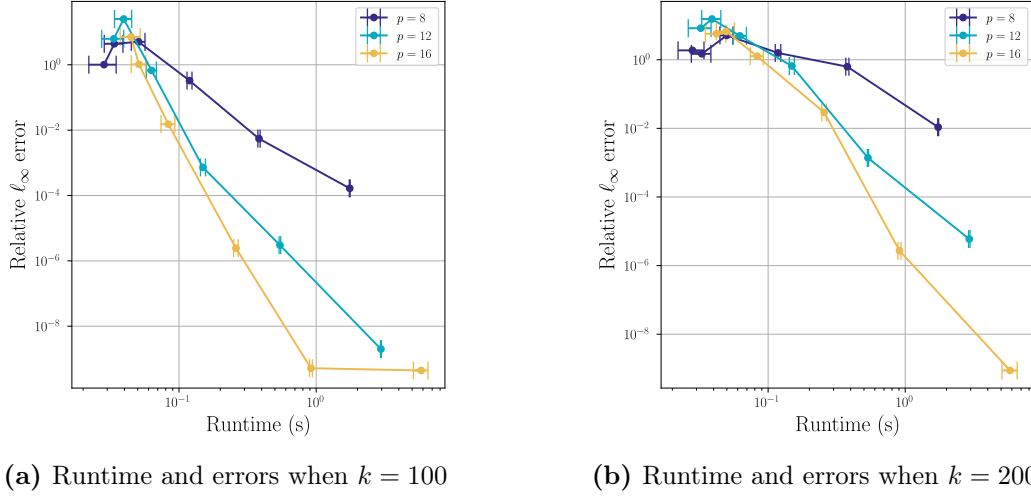
**Figure 5:** Visualizing the solutions of forward scattering problems for the single Gaussian bump scattering potential. Figures 5b and 5c show the real part of the total wave field.



**Figure 6:** Error-runtime study on the single Gaussian bump scattering potential. Using GPU acceleration, our method can rapidly converge to high-accuracy solutions in high-frequency wave scattering problems. The runtime measurements include the runtime of the entire HPS algorithm and the setup and solution of the boundary integral equation enforcing the radiation condition. Horizontal error bars show  $\pm 1$  standard error computed over five trials.



**Figure 7:** Visualizing the solutions of forward scattering problems for the sum of randomly placed Gaussian bumps scattering potential. Figures 7b and 7c show the real part of the total wave field.



**Figure 8:** Error-runtime study on the sum of randomly placed Gaussian bumps scattering potential.

### 5.3 Solving inverse scattering problems with automatic differentiation

Our solver is compatible with the JAX automatic differentiation framework, which allows us to very easily implement gradient-based optimization algorithms using our accelerated HPS solver as a forward model. We consider an inverse scattering task to recover an inhomogeneous scattering potential  $q_\theta$  specified by basis coefficients  $\{\theta_j\}$ :

$$q_\theta(x) = \sum_{b_j \in B_\gamma} \theta_j b_j(x), \quad (11)$$

$$B_\gamma = \left\{ \sin\left(m_1 \frac{\pi}{2} (x_1 - 1)\right) \sin\left(m_2 \frac{\pi}{2} (x_2 - 1)\right) \mid \sqrt{m_1^2 + m_2^2} \leq \gamma \right\}. \quad (12)$$

The basis  $B_\gamma$  is used in [Borges et al. \(2017\)](#) because it spans smooth, bandlimited functions which vanish on the boundary of  $\Omega$ . The goal of this inverse scattering task is to estimate  $\theta$ . Equation (9) describes how  $q_\theta(x)$  affects the scattered wave field  $u_\theta(x)$ . Our forward model  $\mathcal{F} \circ \mathcal{B}$  is a composition of maps where  $\mathcal{B} : \theta \mapsto q_\theta$  is a discrete sine transform and  $\mathcal{F} : q_\theta \mapsto u_\theta$  evaluates the solution of Equation (9) at points  $x^{(j)}$  away from the support of the scattering potential. We choose a ground-truth  $\theta^*$  to be the basis coefficients of the sum of Gaussian bumps scattering potential Equation (10) projected onto  $B_\gamma$  and generate data  $\mathcal{F}[\theta^*]$ . Given this data and our knowledge of the forward model, we wish to recover an estimate of  $\theta^*$ . We can phrase this problem in an optimization framework:

$$\operatorname{argmin}_{\theta \in \mathbb{R}^{N_\theta}} \|(\mathcal{F} \circ \mathcal{B})(\theta) - (\mathcal{F} \circ \mathcal{B})(\theta^*)\|_2^2. \quad (13)$$

We evaluate  $\mathcal{F}$  using our GPU-accelerated fast direct solver. Because our solver is compatible with automatic differentiation, we can solve this optimization problem using gradient-based methods. We first evaluate the accuracy of JAX automatic differentiation by comparing the outputs with the action of  $J[\theta]$ , the Fréchet derivative of  $\mathcal{F} \circ \mathcal{B}$  centered at  $\theta$ . [Borges et al. \(2017\)](#) characterize this derivative in terms of the solution of linear elliptic PDEs which we can compute with our HPS method; we re-state these results in Section D.1 for completeness. Figure 9a shows the convergence of the outputs of automatic differentiation to the Fréchet derivative when holding polynomial order  $p = 16$  fixed and varying the depth of the quadtree,  $L = 1, \dots, 5$ . For leaf size  $h$ , we observe high-order convergence at rate  $O(h^{p-2})$ . For brevity, we defer the details

of this experiment to Section D.2. We remark that high-order accuracy is possible because  $B_\gamma$  defines a smooth global basis which can be represented using the HPS discretization. This is in contrast to differentiating with respect to the values of  $q$  at the HPS discretization points, which effectively introduces nonsmooth coefficient functions. We also note that the use of standard automatic differentiation software requires the presence of the entire computational graph in memory, which means it can not be used in conjunction with the recomputation strategies introduced in Section 4.1.

To solve Equation (13), we propose to use Gauss–Newton iterations for nonlinear least squares problems. This algorithm requires access to  $J[\theta]$ , which we implement using JAX automatic differentiation applied to our HPS solver. We use automatic differentiation to implement the actions  $J[\theta]^*f$  and  $J[\theta]v$  for arbitrary vectors  $f, v$  and estimates  $\theta$ . We pair these subroutines with a sparse linear algebra least-squares solver (Paige and Saunders, 1982) from the SciPy library (Virtanen et al., 2020) to implement the Gauss–Newton algorithm presented in Algorithm 6.

---

**Algorithm 6:** Gauss–Newton iterations for nonlinear least squares problems

---

**Input:** Data  $(\mathcal{F} \circ \mathcal{B})(\theta^*)$ ; Initial estimate  $\theta_0$

- 1  $t \leftarrow 0$
- 2 **while** not converged **do**
- 3   Compose automatic differentiation and our fast direct solver to compile the function  
 $f \mapsto J[\theta_t]^*f$
- 4   Compose automatic differentiation and our fast direct solver to compile the function  
 $v \mapsto J[\theta_t]v$
- 5   Define a linear operator  $J[\theta_t]$  using subroutines  $J[\theta_t]^*f$  and  $J[\theta_t]v$
- 6   Use a least-squares solver to compute  
 $\delta \leftarrow \underset{\delta}{\operatorname{argmin}} \ \|(\mathcal{F} \circ \mathcal{B})(\theta^*) - ((\mathcal{F} \circ \mathcal{B})(\theta_t) + J[\theta_t]\delta)\|_2^2$
- 7    $\theta_{t+1} \leftarrow \theta_t + \delta$
- 8    $t \leftarrow t + 1$

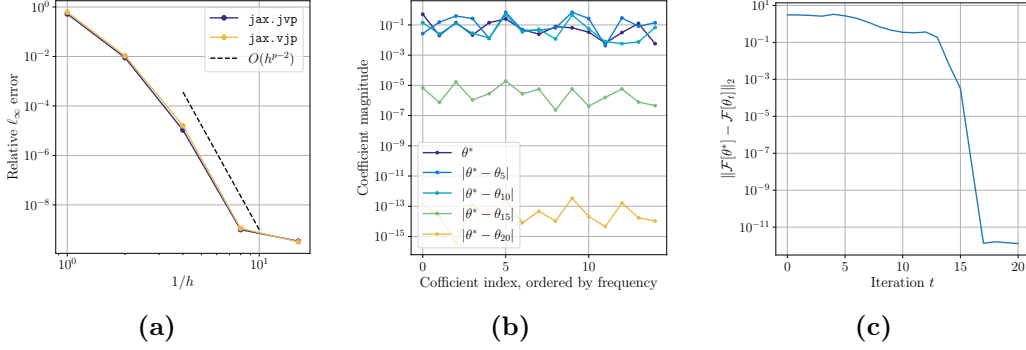
**Result:** Final estimate  $\theta_t$

---

To solve the optimization problem, we use  $\mathcal{B}_\gamma$ ,  $\gamma = 5$ , which gives us  $N_\theta = 15$  optimization variables. Following Borges et al. (2017), we initialize the optimization variables corresponding to the lowest three frequency components to the ground-truth  $\theta_0 = \theta_*$ , and we initialize the other variables at  $\mathbf{0}$ . We run 21 iterations of the Gauss–Newton algorithm. Figure 9 shows the optimization variables take some iterations to approach  $\theta^*$  and then converge superlinearly to near machine precision. Calculating each Gauss–Newton update is fast because of the GPU acceleration of the forward model,  $J[\theta_t]^*f$ , and  $J[\theta_t]v$ . The entire experiment runs in 75 seconds using one H100 GPU.

## 6 Numerical examples in three dimensions

In the three-dimensional examples, we focus on problems with localized regions of high variation. Our adaptive discretization method described in Section 4 is designed for such problems. In these experiments, we use a single Nvidia H100 GPU with 80GB of on-device memory and 200GB of host RAM.



**Figure 9:** Our GPU-accelerated PDE solver can interface with automatic differentiation to rapidly solve inverse problems. In Figure 9a, we show that automatic differentiation applied to our HPS method implementing  $\mathcal{F} \circ \mathcal{B}$  converges at high order to the Fréchet derivative of this operator. Experimental details are available in Section D.2. In Figure 9b, we show the magnitude of the ground-truth coefficients  $\theta^*$  as well as the component-wise errors of intermediate estimates. Figure 9c shows the objective value of the optimization problem, which reaches near machine precision in 17 iterations.

### 6.1 Adaptive refinement on a problem with known solution

In this example, we study the convergence of our method on a three-dimensional problem with a known analytical solution. We build a problem that is solved by a “wavefront” located along a three-dimensional curved surface. The problem is given by

$$\begin{cases} \Delta u(x) = f(x), & x \in [0, 1]^3, \\ u(x) = g(x), & x \in \partial[0, 1]^3. \end{cases} \quad (14)$$

We manufacture a source term  $f(x)$  so the solution takes the form

$$u(x) = u(x_1, x_2, x_3) = \arctan \left( 10 \sqrt{(x_1 - 0.5)^2 + (x_2 - 0.5)^2 + (x_3 - 0.5)^2} - 0.7 \right) \quad (15)$$

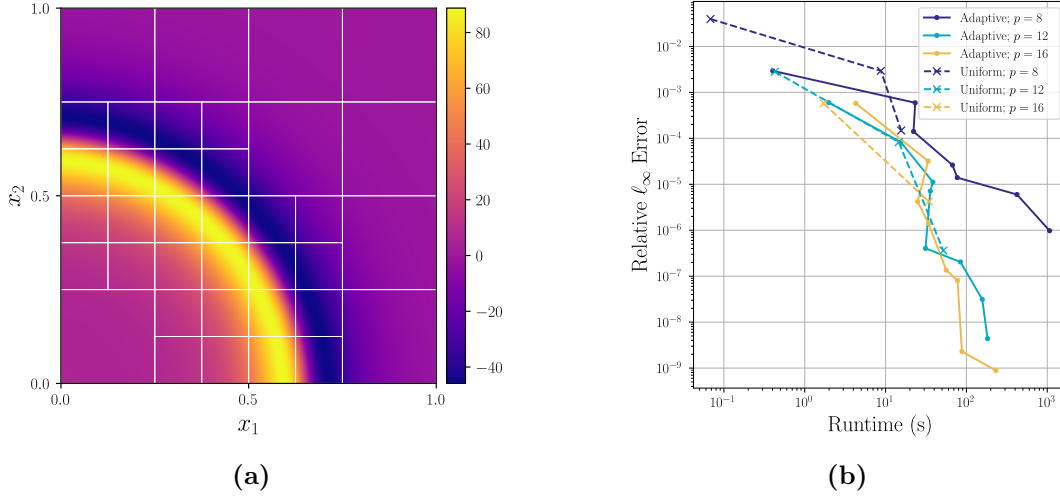
and use samples of this function along the boundary to create our boundary data  $g$ . Figure 10a shows that  $f$  has a localized region of high variation. A uniform discretization strategy cannot adapt to the locality of this problem, but our adaptive discretization strategy can adaptively refine the octree to place a higher density of discretization points in this region. The adaptive discretization also uses larger leaves, where possible, on the parts of the domain with a very smooth source function.

In Figure 10b, we show accuracy versus runtime for both a uniform and adaptive discretization applied to this problem. For all methods, we increased the number of discretization points until saturating the GPU’s memory limit when inverting the highest-level  $\mathbf{D}$  matrix. While the uniform methods are fast on this problem, they are not highly accurate because they cannot use more than  $L = 3$  levels of uniform refinement. Our adaptive method computes solutions with much higher accuracy before saturating the GPU’s memory limit by adaptively placing the discretization points in regions where the source and solution have high variation. Table 1 shows the size of the highest-level merge matrix  $\mathbf{D}$  for selected points on this graph.

### 6.2 Linearized Poisson–Boltzmann equation

An example application where the data and solution have local regions of high variation is the linearized Poisson–Boltzmann equation, a model of the electrostatic properties of a molecule





**Figure 10:** Adaptive refinement allows for more accuracy before encountering memory bottlenecks. Figure 10a shows the source function restricted to the  $x_3 = 0$  plane and the adaptive mesh formed with tolerance  $1 \times 10^{-7}$  and Chebyshev parameter  $p = 16$ . In Figure 10b, we study the runtime and error of different refinement strategies applied to Equation (14). For each step of the uniform refinement curve, we refined the uniform grid by one more level. For the adaptive refinement curve, we decreased the adaptive refinement tolerance by a factor of 10. For all methods, we refined until running out of memory on the GPU during the build stage.

in a solution. This can, for example, be used to compute the stability of a given molecular configuration in a solution. A standard model, developed in Grant et al. (2001), starts with atoms represented by point charges  $\{z^{(i)}\}_{i=1}^{N_z}$ ,  $z^{(i)} \in \mathbb{R}^3$ . These atoms give rise to a charge distribution  $\rho(x)$ ,

$$\rho(x) = \sum_{i=1}^{N_z} e^{-\delta \|x - z^{(i)}\|^2}, \quad (16)$$

and a spatially-varying permittivity function  $\varepsilon(x)$ ,

$$\varepsilon(x) = \epsilon_0 + (\epsilon_\infty - \epsilon_0)e^{-A\rho(x)}. \quad (17)$$

We use parameters  $N_z = 50$ ,  $\delta = 45$ ,  $\epsilon_0 = 16$ ,  $\epsilon_\infty = 100$ , and  $A = 10$  (Grant et al., 2001). The atomic centers  $\{z^{(i)}\}_{i=1}^{N_z}$  are drawn uniformly from the box  $[-0.5, 0.5]^3$ . We can now model the electrostatic potential  $u(x)$  which is implicitly defined by the linearized Poisson–Boltzmann equation:

$$\begin{cases} \nabla \cdot (\varepsilon(x) \cdot \nabla u(x)) = -\rho(x), & x \in [-1, 1]^3, \\ u(x) = 0, & x \in \partial[-1, 1]^3. \end{cases} \quad (18)$$

Existing approaches, such as a fast integral equation method (Vico et al., 2016) and finite difference schemes (Nicholls and Honig, 1991; Colmenares et al., 2014) solve a simplified version of Equation (18), where the permittivity function is replaced by one derived from van der Waals surfaces:

$$\varepsilon_{\text{vdW}}(x) = q(x)\epsilon_0 + (1 - q(x))(\epsilon_\infty - \epsilon_0), \quad (19)$$

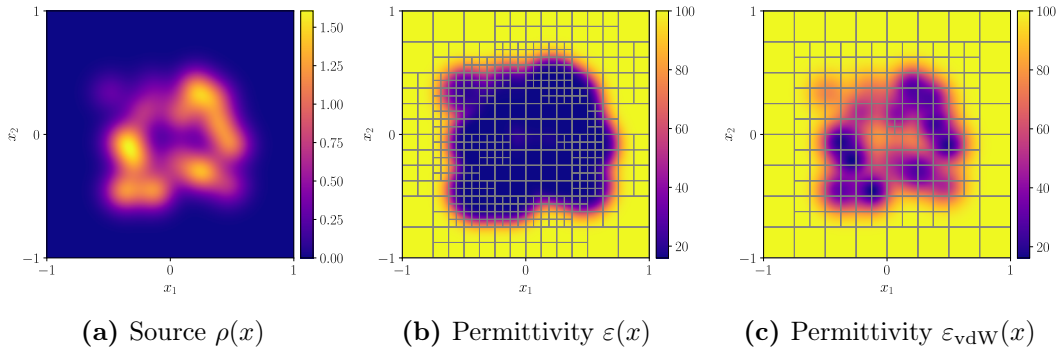
$$q(x) = 1 - \prod_{i=1}^{N_z} \left[ 1 - e^{-\delta \|x - z^{(i)}\|^2} \right]. \quad (20)$$

$p$	Tolerance	$n_{\text{leaves}}$	$N$	Max Depth	Runtime (s)	$p$	Tolerance	$n_{\text{leaves}}$	$N$	Max Depth	Runtime (s)
8	$10^{-1}$	358	183,296	3	48.4	8	$1 \times 10^{-3}$	1,093	559,616	4	154.4
8	$10^{-2}$	1,436	735,232	4	177.8	8	$1 \times 10^{-4}$	1,737	889,344	5	227.2
8	$10^{-3}$	1,884	964,608	5	218.9	8	$1 \times 10^{-5}$	5,111	2,616,832	5	869.0
8	$10^{-4}$	7,659	3,921,408	5	1,523.0	8	$1 \times 10^{-6}$	9,423	4,824,576	5	OOM
10	$10^{-1}$	253	253,000	3	52.5	10	$1 \times 10^{-3}$	435	435,000	4	94.5
10	$10^{-2}$	449	449,000	4	76.3	10	$1 \times 10^{-4}$	1,016	1,016,000	4	193.3
10	$10^{-3}$	1,625	1,625,000	4	243.4	10	$1 \times 10^{-5}$	1,485	1,485,000	4	241.7
10	$10^{-4}$	1,982	1,982,000	5	319.1	10	$1 \times 10^{-6}$	2,605	2,605,000	5	483.0
12	$10^{-1}$	99	171,072	3	41.5	12	$1 \times 10^{-3}$	274	473,472	3	93.7
12	$10^{-2}$	386	667,008	3	91.0	12	$1 \times 10^{-4}$	477	824,256	4	145.8
12	$10^{-3}$	715	1,235,520	4	190.6	12	$1 \times 10^{-5}$	974	1,683,072	4	275.2
12	$10^{-4}$	1,695	2,928,960	4	OOM	12	$1 \times 10^{-6}$	1,366	2,360,448	4	OOM
16	$10^{-1}$	64	262,144	2	53.9	16	$1 \times 10^{-3}$	127	520,192	3	119.7
16	$10^{-2}$	204	835,584	3	136.2	16	$1 \times 10^{-4}$	239	978,944	3	172.2
16	$10^{-3}$	365	1,495,040	3	OOM	16	$1 \times 10^{-5}$	323	1,323,008	3	227.0
16	$10^{-4}$	470	1,925,120	4	OOM	16	$1 \times 10^{-6}$	456	1,867,776	4	OOM

**Table 2:** Resource usage statistics for using our 3D adaptive HPS method applied to the linearized Poisson–Boltzmann equation (Equation (18)) using permittivity  $\varepsilon(x)$  (left) and simplified permittivity  $\varepsilon_{\text{vdW}}(x)$  (right). We use “OOM” to indicate which discretizations caused out of memory errors when inverting the final  $\mathbf{D}$  matrix.

$\varepsilon_{\text{vdW}}(x)$  is an easier function to resolve to high accuracy as it lacks the steep gradients observed in  $\varepsilon(x)$ . However, Grant et al. (2001) reports “experimentation with a number of dielectric mapping functions using  $[\varepsilon_{\text{vdW}}]$  produced dielectric functions that increase toward solvent values far too rapidly with distance from atomic centers,” and “ $[\varepsilon_{\text{vdW}}]$  also produced undesired patches of high dielectric inside proteins.” We use our GPU-accelerated adaptive HPS method to solve Equation (18) with both permittivity models.

To form an adaptive discretization for this problem, we refine a discretization tree given the charge distribution  $\rho$ , the permittivity  $\varepsilon$ , and the components of  $\nabla\varepsilon$ . Figure 11 shows a 2D slice of the charge distribution and permittivity, along with the discretization found by this refinement process. Table 2 gives statistics about the discretization and runtime for a range of tolerances and Chebyshev parameters  $p$ . In this table,  $n_{\text{leaves}}$  is the number of leaves of the resulting discretization tree,  $N = n_{\text{leaves}}p^3$  is the total number of discretization points, Max Depth is the maximum number of levels of refinement in the discretization tree, and Runtime measures the wall-clock time in seconds to compute the adaptive discretization and execute all parts of the HPS algorithm.



**Figure 11:** Visualizing the variable coefficients and source term of our problem. These plots show the source function  $\rho(x)$  and permittivity functions  $\varepsilon(x)$  and  $\varepsilon_{\text{vdW}}(x)$  restricted to the plane  $x_3 = 0$ . The adaptive discretizations formed using  $p = 8$  and tolerance  $1 \times 10^{-4}$  are shown. This adaptive discretization is found by forming the union of meshes adaptively refined on the source, the permittivity, and components of the gradient of the permittivity.

## 7 Conclusion

This paper presents methods for efficiently accelerating HPS algorithms using general-purpose GPUs. Because there is a large amount of inherent parallelism in the structure of the HPS algorithms, they are a natural target for GPU acceleration once adjustments are made to reduce memory complexity. We introduce methods for reducing the memory footprint of HPS algorithms for problems in two and three dimensions.

This work leaves open important questions and avenues for improvement. While our method can efficiently interface with automatic differentiation, we could, in principle, gain much more efficiency by implementing custom automatic differentiation rules to reuse the precomputed solution operators, like those derived in [Borges et al. \(2017\)](#). Our methods of reducing memory complexity could be pushed further by using a hybrid approach, i.e., by performing a few levels of merging via dense linear algebra and then relying on a sparse direct solver such as [Yesypenko and Martinsson \(2024b\)](#) or [Kump et al. \(2025\)](#) for higher-level merges. We hypothesize such a hybrid approach would greatly reduce runtimes for very large problems by reducing the ranks needed to accurately resolve the sparse system matrix. Finally, our methods could be extended to unstructured meshes and surfaces, a setting where nonuniform merge operations make parallel implementations challenging.

## Author Contributions

- **Owen Melia:** Conceptualization; Methodology; Software; Writing - original draft.
- **Daniel Fortunato:** Conceptualization; Methodology; Supervision; Writing - review & editing.
- **Jeremy Hoskins:** Conceptualization; Methodology; Supervision; Writing - review & editing.
- **Rebecca Willett:** Conceptualization; Methodology; Supervision; Funding acquisition; Project administration; Writing - review & editing.

## Data Availability

Our open-source JAX implementation is publicly available at <https://github.com/meliao/jaxhps>.

## Acknowledgment

The authors would like to thank Manas Rachh, Leslie Greengard, and Olivia Tsang for many useful discussions. The authors are grateful to the Flatiron Institute for providing the computational resources used to conduct the experiments in this work. OM and RW gratefully acknowledge the support of NSF DMS-2023109, DOE DE-SC0022232, the Physics Frontier Center for Living Systems funded by the National Science Foundation (PHY-2317138), and the support of the Margot and Tom Pritzker Foundation. The Flatiron Institute is a division of the Simons Foundation.

## References

- Abdelfattah, A., Barra, V., Beams, N., Bleile, R., Brown, J., Camier, J.S., Carson, R., Chalmers, N., Dobrev, V., Dudouit, Y., Fischer, P., Karakus, A., Kerkemeier, S., Kolev, T., Lan, Y.H., Merzari, E., Min, M., Phillips, M., Rathnayake, T., Rieben, R., Stitt, T., Tomboulides, A., Tomov, S., Tomov, V., Vargas, A., Warburton, T., Weiss, K., 2021. GPU algorithms for efficient exascale discretizations. *Parallel Computing* 108, 102841. doi:doi:10.1016/j.parco.2021.102841.
- Abdelfattah, A., Ghysels, P., Boukaram, W., Tomov, S., Li, X.S., Dongarra, J., 2022. Addressing irregular patterns of matrix computations on GPUs and their impact on applications powered by sparse direct solvers, in: SC22: International Conference for High Performance Computing, Networking, Storage and Analysis, IEEE, Dallas, TX, USA. pp. 1–14. URL: <https://ieeexplore.ieee.org/document/10046092/>, doi:doi:10.1109/SC41404.2022.00031.
- Ansel, J., Yang, E., He, H., Gimelshein, N., Jain, A., Voznesensky, M., Bao, B., Bell, P., Berard, D., Burovski, E., Chauhan, G., Chourdia, A., Constable, W., Desmaison, A., DeVito, Z., Ellison, E., Feng, W., Gong, J., Gschwind, M., Hirsh, B., Huang, S., Kalambarkar, K., Kirsch, L., Lazos, M., Lezcano, M., Liang, Y., Liang, J., Lu, Y., Luk, C., Maher, B., Pan, Y., Puhersch, C., Reso, M., Saroufim, M., Siraichi, M.Y., Suk, H., Suo, M., Tillet, P., Wang, E., Wang, X., Wen, W., Zhang, S., Zhao, X., Zhou, K., Zou, R., Mathews, A., Chanan, G., Wu, P., Chintala, S., 2024. PyTorch 2: Faster machine learning through dynamic Python bytecode transformation and graph compilation. URL: <https://pytorch.org/assets/pytorch2-2.pdf>, doi:doi:10.1145/3620665.3640366. 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2 (ASPLOS '24).
- Askham, T., Rachh, M., O'Neil, M., Hoskins, J., Fortunato, D., Jiang, S., Fryklund, F., Goodwill, T., Wang, H.Y., Zhu, H., 2024. chunkIE: A MATLAB integral equation toolbox. URL: <https://github.com/fastalgorithms/chunkie>.
- Beams, N.N., Gillman, A., Hewett, R.J., 2020. A parallel shared-memory implementation of a high-order accurate solution technique for variable coefficient helmholtz problems. *Computers & Mathematics with Applications* 79, 996–1011. doi:doi:10.1016/j.camwa.2019.08.019.
- Beck, T., Canzani, Y., Marzuola, J.L., 2022. Quantitative bounds on impedance-to-impedance operators with applications to fast direct solvers for PDEs. *Pure and Applied Analysis* 4, 225–256. URL: <https://msp.org/paa/2022/4-2/p02.xhtml>, doi:doi:10.2140/paa.2022.4.225. publisher: Mathematical Sciences Publishers.
- Borges, C., Gillman, A., Greengard, L., 2017. High resolution inverse scattering in two dimensions using recursive linearization. *SIAM Journal on Imaging Sciences* 10, 641–664. doi:doi:10.1137/16M1093562.
- Bradbury, J., Frostig, R., Hawkins, P., Johnson, M.J., Leary, C., Maclaurin, D., Necula, G., Paszke, A., VanderPlas, J., Wanderman-Milne, S., Zhang, Q., 2018. JAX: Composable transformations of Python+NumPy programs. URL: <http://github.com/jax-ml/jax>.
- Chipman, D., Calhoun, D., Burstedde, C., 2024. A fast direct solver for elliptic PDEs on a hierarchy of adaptively refined quadrees. URL: <http://arxiv.org/abs/2402.14936>. arXiv:2402.14936 [cs, math].
- Colmenares, J., Ortiz, J., Rocchia, W., 2014. GPU linear and non-linear Poisson–Boltzmann solver module for DelPhi. *Bioinformatics* 30, 569–570. URL: <https://doi.org/10.1093/bioinformatics/btt699>, doi:doi:10.1093/bioinformatics/btt699.

- Dao, T., Fu, D., Ermon, S., Rudra, A., Ré, C., 2022. FlashAttention: Fast and memory-efficient exact attention with IO-awareness, in: Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., Oh, A. (Eds.), *Advances in Neural Information Processing Systems*, Curran Associates, Inc.. p. 16344–16359. URL: [https://proceedings.neurips.cc/paper\\_files/paper/2022/file/67d57c32e20fd0a7a302cb81d36e40d5-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/67d57c32e20fd0a7a302cb81d36e40d5-Paper-Conference.pdf).
- Diao, K., Li, Z., Grumitt, R.D.P., Mao, Y., 2024. *synax*: A differentiable and gpu-accelerated synchrotron simulation package URL: <http://arxiv.org/abs/2410.01136>, doi:doi:10.48550/arXiv.2410.01136. arXiv:2410.01136 [astro-ph].
- Ernst, O.G., Gander, M.J., 2012. Why it is Difficult to Solve Helmholtz Problems with Classical Iterative Methods, in: Graham, I.G., Hou, T.Y., Lakkis, O., Scheichl, R. (Eds.), *Numerical Analysis of Multiscale Problems*. Springer, Berlin, Heidelberg, pp. 325–363. URL: [https://doi.org/10.1007/978-3-642-22061-6\\_10](https://doi.org/10.1007/978-3-642-22061-6_10), doi:doi:10.1007/978-3-642-22061-6\_10.
- Fortunato, D., 2024. A high-order fast direct solver for surface PDEs. *SIAM Journal on Scientific Computing* 46, A2582–A2606. URL: <https://epubs.siam.org/doi/abs/10.1137/22M1525259>, doi:doi:10.1137/22M1525259. publisher: Society for Industrial and Applied Mathematics.
- Fortunato, D., Hale, N., Townsend, A., 2021. The ultraspherical spectral element method. *Journal of Computational Physics* 436, 110087. doi:doi:10.1016/j.jcp.2020.110087.
- Geldermans, P., Gillman, A., 2019. An adaptive high order direct solution technique for elliptic boundary value problems. *SIAM Journal on Scientific Computing* 41, A292–A315. URL: <https://epubs.siam.org/doi/abs/10.1137/17M1156320>, doi:doi:10.1137/17M1156320. publisher: Society for Industrial and Applied Mathematics.
- George, A., 1973. Nested dissection of a regular finite element mesh. *SIAM Journal on Numerical Analysis* 10, 345–363. URL: <https://epubs.siam.org/doi/10.1137/0710032>, doi:doi:10.1137/0710032. publisher: Society for Industrial and Applied Mathematics.
- Georgescu, S., Chow, P., Okuda, H., 2013. GPU acceleration for FEM-based structural analysis. *Archives of Computational Methods in Engineering* 20, 111–121. URL: <https://doi.org/10.1007/s11831-013-9082-8>, doi:doi:10.1007/s11831-013-9082-8.
- Ghysels, P., Synk, R., 2022. High performance sparse multifrontal solvers on modern GPUs. *Parallel Computing* 110, 102897. URL: <https://www.sciencedirect.com/science/article/pii/S0167819122000059>, doi:doi:10.1016/j.parco.2022.102897.
- Gillman, A., Barnett, A.H., Martinsson, P.G., 2015. A spectrally accurate direct solution technique for frequency-domain scattering problems with variable media. *BIT Numerical Mathematics* 55, 141–170. URL: <https://doi.org/10.1007/s10543-014-0499-8>, doi:doi:10.1007/s10543-014-0499-8.
- Gillman, A., Martinsson, P.G., 2014. A direct solver with  $O(N)$  complexity for variable coefficient elliptic PDEs discretized via a high-order composite spectral collocation method. *SIAM Journal on Scientific Computing* 36, A2023–A2046. URL: <https://epubs.siam.org/doi/10.1137/130918988>, doi:doi:10.1137/130918988. publisher: Society for Industrial and Applied Mathematics.
- Grant, J.A., Pickup, B.T., Nicholls, A., 2001. A smooth permittivity function for Poisson–Boltzmann solvation methods. *Journal of Computational Chemistry* 22, 608–640. URL: <https://onlinelibrary.wiley.com/doi/abs/10.1002/jcc.1032>, doi:doi:10.1002/jcc.1032.

- Gu, A., Dao, T., 2024. Mamba: Linear-time sequence modeling with selective state spaces. URL: <https://openreview.net/forum?id=tEYskw1VY2#discussion>.
- Hao, S., Martinsson, P.G., 2016. A direct solver for elliptic PDEs in three dimensions based on hierarchical merging of Poincaré–Steklov operators. *Journal of Computational and Applied Mathematics* 308, 419–434. URL: <https://www.sciencedirect.com/science/article/pii/S0377042716302308>, doi:doi:10.1016/j.cam.2016.05.013.
- Ho, K.L., Greengard, L., 2012. A fast direct solver for structured linear systems by recursive skeletonization. *SIAM Journal on Scientific Computing* 34, A2507–A2532. URL: <https://epubs.siam.org/doi/10.1137/120866683>, doi:doi:10.1137/120866683. publisher: Society for Industrial and Applied Mathematics.
- Kidger, P., 2021. On Neural Differential Equations. PhD Thesis. University of Oxford.
- Kolev, T., Fischer, P., Min, M., Dongarra, J., Brown, J., Dobrev, V., Warburton, T., Tomov, S., Shephard, M.S., Abdelfattah, A., Barra, V., Beams, N., Camier, J.S., Chalmers, N., Dudouit, Y., Karakus, A., Karlin, I., Kerkemeier, S., Lan, Y.H., Medina, D., Merzari, E., Obabko, A., Pazner, W., Rathnayake, T., Smith, C.W., Spies, L., Swirydowicz, K., Thompson, J., Tomboulides, A., Tomov, V., 2021. Efficient exascale discretizations: High-order finite element methods. *The International Journal of High Performance Computing Applications* 35, 527–552. doi:doi:10.1177/10943420211020803.
- Kopriva, D.A., 1998. A staggered-grid multidomain spectral method for the compressible Navier–Stokes equations. *Journal of Computational Physics* 143, 125–158. URL: <https://www.sciencedirect.com/science/article/pii/S0021999198959563>, doi:doi:10.1006/jcph.1998.5956.
- Kump, J., Yesypenko, A., Martinsson, P.G., 2025. A two-level direct solver for the hierarchical Poincaré–Steklov method URL: <http://arxiv.org/abs/2503.04033>, doi:doi:10.48550/arXiv.2503.04033. arXiv:2503.04033 [math].
- Leary, C., Wang, T., 2017. XLA: TensorFlow, Compiled!
- Li, X.S., Demmel, J.W., 2003. SuperLU\_DIST: A scalable distributed-memory sparse direct solver for unsymmetric linear systems. *ACM Trans. Math. Softw.* 29, 110–140. URL: <https://dl.acm.org/doi/10.1145/779359.779361>, doi:doi:10.1145/779359.779361.
- Lucero Lorca, J.P., Beams, N., Beecroft, D., Gillman, A., 2024. An iterative solver for the HPS discretization applied to three dimensional Helmholtz problems. *SIAM Journal on Scientific Computing* 46, A80–A104. URL: <https://epubs.siam.org/doi/full/10.1137/21M1463380>, doi:doi:10.1137/21M1463380. publisher: Society for Industrial and Applied Mathematics.
- Martinsson, P.G., 2013. A direct solver for variable coefficient elliptic PDEs discretized via a composite spectral collocation method. *Journal of Computational Physics* 242, 460–479. URL: <https://www.sciencedirect.com/science/article/pii/S0021999113001320>, doi:doi:10.1016/j.jcp.2013.02.019.
- Martinsson, P.G., 2015. The Hierarchical Poincaré–Steklov (HPS) solver for elliptic PDEs: A tutorial. URL: <http://arxiv.org/abs/1506.01308>. arXiv:1506.01308 [math].
- Martinsson, P.G., 2019. Fast Direct Solvers for Elliptic PDEs. Society for Industrial and Applied Mathematics, Philadelphia, PA. URL: <https://epubs.siam.org/doi/book/10.1137/1.9781611976045>, doi:doi:10.1137/1.9781611976045.



- Nicholls, A., Honig, B., 1991. A rapid finite difference algorithm, utilizing successive over-relaxation to solve the Poisson–Boltzmann equation. *Journal of Computational Chemistry* 12, 435–445. URL: <https://onlinelibrary.wiley.com/doi/abs/10.1002/jcc.540120405>, doi:doi:10.1002/jcc.540120405.
- Paige, C.C., Saunders, M.A., 1982. LSQR: An algorithm for sparse linear equations and sparse least squares. *ACM Transactions on Mathematical Software* 8, 43–71. URL: <https://dl.acm.org/doi/10.1145/355984.355989>, doi:doi:10.1145/355984.355989.
- Pfeiffer, H.P., Kidder, L.E., Scheel, M.A., Teukolsky, S.A., 2003. A multidomain spectral method for solving elliptic equations. *Computer Physics Communications* 152, 253–273. URL: <https://www.sciencedirect.com/science/article/pii/S0010465502008470>, doi:doi:10.1016/S0010-4655(02)00847-0.
- Trefethen, L.N., 2000. *Spectral Methods in MATLAB. Software, Environments, and Tools*, Society for Industrial and Applied Mathematics. URL: <https://epubs-siam-org.proxy.uchicago.edu/doi/book/10.1137/1.9780898719598>, doi:doi:10.1137/1.9780898719598.
- Vico, F., Greengard, L., Ferrando, M., 2016. Fast convolution with free-space Green’s functions. *Journal of Computational Physics* 323, 191–203. URL: <https://www.sciencedirect.com/science/article/pii/S0021999116303230>, doi:doi:10.1016/j.jcp.2016.07.028.
- Virtanen, P., Gommers, R., Oliphant, T.E., Haberland, M., Reddy, T., Cournapeau, D., Burovski, E., Peterson, P., Weckesser, W., Bright, J., van der Walt, S.J., Brett, M., Wilson, J., Millman, K.J., Mayorov, N., Nelson, A.R.J., Jones, E., Kern, R., Larson, E., Carey, C.J., Polat, İ., Feng, Y., Moore, E.W., VanderPlas, J., Laxalde, D., Perktold, J., Cimrman, R., Henriksen, I., Quintero, E.A., Harris, C.R., Archibald, A.M., Ribeiro, A.H., Pedregosa, F., van Mulbregt, P., SciPy 1.0 Contributors, 2020. SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nature Methods* 17, 261–272. doi:doi:10.1038/s41592-019-0686-2.
- Xue, T., Liao, S., Gan, Z., Park, C., Xie, X., Liu, W.K., Cao, J., 2023. JAX-FEM: A differentiable GPU-accelerated 3D finite element solver for automatic inverse design and mechanistic data science. *Computer Physics Communications* 291, 108802. URL: <https://www.sciencedirect.com/science/article/pii/S0010465523001479>, doi:doi:10.1016/j.cpc.2023.108802.
- Yang, B., Hesthaven, J.S., 2000. Multidomain pseudospectral computation of Maxwell’s equations in 3-D general curvilinear coordinates. *Applied Numerical Mathematics* 33, 281–289. URL: <https://www.sciencedirect.com/science/article/pii/S016892749900094X>, doi:doi:10.1016/S0168-9274(99)00094-X.
- Yesypenko, A., Martinsson, P.G., 2024a. GPU optimizations for the hierarchical Poincaré–Steklov Scheme, in: Dostál, Z., Kozubek, T., Klawonn, A., Langer, U., Pavarino, L.F., Šístek, J., Widlund, O.B. (Eds.), *Domain Decomposition Methods in Science and Engineering XXVII*, Springer Nature Switzerland, Cham. pp. 519–528. doi:doi:10.1007/978-3-031-50769-4\_62.
- Yesypenko, A., Martinsson, P.G., 2024b. SlabLU: a two-level sparse direct solver for elliptic PDEs. *Advances in Computational Mathematics* 50, 90. URL: <https://doi.org/10.1007/s10444-024-10176-x>, doi:doi:10.1007/s10444-024-10176-x.

## A Full algorithms for 2D problems using DtN matrices

In this section, we describe the details for the two-dimensional version of our method which merges DtN matrices. In this version of the algorithm, the outgoing boundary data  $\mathbf{h}$  tabulates the outward-pointing normal derivative of the particular solution, and the incoming boundary data  $\mathbf{g}$  tabulates the homogenous solution values restricted to patch boundaries.  $\mathbf{T}$  is a Dirichlet-to-Neumann matrix.

In this section, we use  $\mathbf{I}_{a \times a}$  to denote the identity matrix of shape  $a \times a$  and  $\mathbf{0}_d$  to denote a length- $d$  vector filled with 0's. When defining matrices blockwise, we use  $\mathbf{0}$  to denote a block filled with 0's, and assume the shape of the block can be determined from the nonzero blocks sharing the same rows and columns.

### A.1 Local solve stage

Recall from Section 3.1 that we use a tensor product of order- $p$  Chebyshev–Lobatto points to discretize the interior of each leaf. This results in a grid with  $p^2$  discretization points, and  $4p - 4$  of these points lie on the boundary of the leaf. We use order- $q$  Gauss–Legendre points to discretize each side of the leaf's boundary, so there are  $4q$  boundary points in total. Thus, to translate the information between the interior and boundary of each leaf, we need to compute spectral differentiation matrices and matrices interpolating between the  $p$  Chebyshev and  $q$  Gauss points. In particular, we need to precompute the following matrices:

- $\mathbf{P}$ , with shape  $4p - 4 \times 4q$ , is the operator mapping data sampled on the Gauss boundary points to data sampled on the  $4p - 4$  Chebyshev points located on the boundary of the leaf. This matrix is constructed using a barycentric Lagrange interpolation matrix mapping from Gauss to Chebyshev points on one side of the leaf; this interpolation matrix is repeated for the other sides. Rows corresponding to the Chebyshev points on the corners of the leaf average the contribution from the two adjoining panels.
- $\mathbf{Q}$ , with shape  $4q \times p^2$ , performs spectral differentiation on the  $p^2$  Chebyshev points followed by interpolation to the Gauss boundary points. This matrix is formed by stacking the relevant rows of Chebyshev spectral differentiation matrices to form an operator which evaluates normal derivatives on the  $4p - 4$  boundary Chebyshev points, and then composing this differentiation operator with a matrix formed from barycentric Lagrange interpolation matrix blocks. These interpolation matrices each map from one Chebyshev panel to one Gauss panel.

To work with  $\mathbf{L}^{(i)}$ , the discretization of the differential operator on leaf  $i$ , it is useful to identify  $I_i$  and  $I_e$ , the sets of discretization points corresponding to the  $(p - 2)^2$  interior and  $4p - 4$  exterior Chebyshev points, respectively. Now, we can fully describe the local solve stage in Algorithm 7.

### A.2 Merge stage

In the 2D merge stage, we are merging four nodes  $\Omega_a, \Omega_b, \Omega_c$ , and  $\Omega_d$ , which have *exterior* and *interior* discretization points. We label the exterior boundary sections 1, 2, 3, and 4, and we label the interior boundary sections 5, 6, 7, and 8. See Figure 12 for a diagram of the different boundary parts. Because the merge stage operates completely on data discretized using Gauss–Legendre panels, there are no discretization points at the corners of nodes. This means each discretization point belongs to exactly one part of the boundary. During this stage of the algorithm, we will be indexing rows and columns of the Dirichlet-to-Neumann matrices according to these boundary sections. For example, we use  $\mathbf{T}_{1,5}^{(a)}$  to indicate the submatrix of

---

**Algorithm 7:** 2D DtN local solve stage.

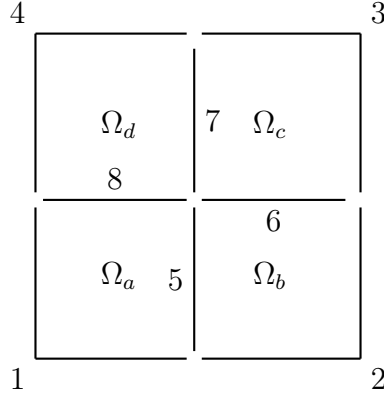
---

**Input:** Discretized differential operators  $\{\mathbf{L}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$ ; discretized source vectors  $\{\mathbf{f}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$ ; precomputed interpolation and differentiation matrices  $\mathbf{P}$  and  $\mathbf{Q}$

- 1 **for**  $i = 1, \dots, n_{\text{leaves}}$  **do**
- 2   Invert  $\mathbf{L}^{(i)}(I_i, I_i)$
- 3    $\mathbf{Y}^{(i)} = \begin{bmatrix} \mathbf{I}_{4p-4 \times 4p-4} \\ -(\mathbf{L}^{(i)}(I_i, I_i))^{-1} \mathbf{L}^{(i)}(I_i, I_e) \end{bmatrix} \mathbf{P}$
- 4    $\mathbf{v}^{(i)} = \begin{bmatrix} \mathbf{0}_{4p-4} \\ -(\mathbf{L}^{(i)}(I_i, I_i))^{-1} \mathbf{f}^{(i)}(I_i) \end{bmatrix}$
- 5    $\mathbf{T}^{(i)} = \mathbf{Q} \mathbf{Y}^{(i)}$
- 6    $\mathbf{h}^{(i)} = \mathbf{Q} \mathbf{v}^{(i)}$

**Result:** Poincaré–Steklov matrices  $\{\mathbf{T}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$ ; outgoing boundary data  $\{\mathbf{h}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$ ; interior solution matrices  $\{\mathbf{Y}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$ , leaf-level particular solutions  $\{\mathbf{v}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$

---



**Figure 12:** Visualizing boundary elements 1 through 8 for two-dimensional merges.

node  $a$ 's DtN matrix which maps from boundary section 5 to boundary section 1. Suppose each side of  $\Omega_a, \Omega_b, \Omega_c$ , and  $\Omega_d$  is discretized with  $n_{\text{side}}$  discretization points; in this case  $\mathbf{T}_{1,5}^{(a)}$  will have shape  $2n_{\text{side}} \times n_{\text{side}}$ .

To implement the merge stage, we use sets of constraints to solve for a mapping from given  $\mathbf{g}_{\text{ext}}$  to unknown  $\mathbf{g}_{\text{int}}$ . These are vectors tabulating the homogeneous solution along the exterior and interior boundary parts. First are constraints specifying that the solution to the PDE is continuous:

$$\mathbf{u}_{\text{ext}} = \mathbf{A} \mathbf{g}_{\text{ext}} + \mathbf{B} \mathbf{g}_{\text{int}} + \mathbf{h}_{\text{ext}}^{(\text{child})}. \quad (21)$$

In this equation  $\mathbf{u}_{\text{ext}}$  is interpreted as the outward-pointing normal derivative of the solution to the PDE restricted to boundary elements 1, 2, 3, and 4 with boundary data specified by  $\mathbf{g}_{\text{ext}}$ .

In this set of constraints, we define:

$$\mathbf{h}_{\text{ext}}^{(\text{child})} = \begin{bmatrix} \mathbf{h}_1^{(a)} \\ \mathbf{h}_2^{(b)} \\ \mathbf{h}_3^{(c)} \\ \mathbf{h}_4^{(d)} \end{bmatrix}, \quad (22)$$

$$\mathbf{A} = \begin{bmatrix} \mathbf{T}_{1,1}^{(a)} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{T}_{2,2}^{(b)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{T}_{3,3}^{(c)} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{T}_{4,4}^{(d)} \end{bmatrix}, \quad (23)$$

$$\mathbf{B} = \begin{bmatrix} \mathbf{T}_{1,5}^{(a)} & \mathbf{0} & \mathbf{0} & \mathbf{T}_{1,8}^{(a)} \\ \mathbf{T}_{2,5}^{(b)} & \mathbf{T}_{2,6}^{(b)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{T}_{3,6}^{(c)} & \mathbf{T}_{3,7}^{(c)} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{T}_{4,7}^{(d)} & \mathbf{T}_{4,8}^{(d)} \end{bmatrix}. \quad (24)$$

A second set of constraints enforces that the outward-pointing normal derivatives from neighboring nodes should sum to zero, which is equivalent to enforcing the continuity of the first derivative along the merge interfaces in their respective Cartesian directions. To that end, we use constraints  $\mathbf{u}_5^{(a)} + \mathbf{u}_5^{(b)} = \mathbf{0}_{n_{\text{side}}}$ ,  $\mathbf{u}_6^{(b)} + \mathbf{u}_6^{(c)} = \mathbf{0}_{n_{\text{side}}}$ , and so on. This gives us an equation:

$$\mathbf{0}_{4n_{\text{side}}} = \mathbf{C}\mathbf{g}_{\text{ext}} + \mathbf{D}\mathbf{g}_{\text{int}} + \mathbf{h}_{\text{int}}^{(\text{child})}. \quad (25)$$

In Equation (25), we define

$$\mathbf{h}_{\text{int}}^{(\text{child})} = \begin{bmatrix} \mathbf{h}_5^{(a)} + \mathbf{h}_5^{(b)} \\ \mathbf{h}_6^{(b)} + \mathbf{h}_6^{(c)} \\ \mathbf{h}_7^{(c)} + \mathbf{h}_7^{(d)} \\ \mathbf{h}_8^{(d)} + \mathbf{h}_8^{(a)} \end{bmatrix}, \quad (26)$$

$$\mathbf{C} = \begin{bmatrix} \mathbf{T}_{5,1}^{(a)} & \mathbf{T}_{5,2}^{(b)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{T}_{6,2}^{(b)} & \mathbf{T}_{6,3}^{(c)} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{T}_{7,3}^{(c)} & \mathbf{T}_{7,4}^{(d)} \\ \mathbf{T}_{8,1}^{(a)} & \mathbf{0} & \mathbf{0} & \mathbf{T}_{8,4}^{(d)} \end{bmatrix}, \quad (27)$$

$$\mathbf{D} = \begin{bmatrix} \mathbf{T}_{5,5}^{(a)} + \mathbf{T}_{5,5}^{(b)} & \mathbf{T}_{5,6}^{(b)} & \mathbf{0} & \mathbf{T}_{5,8}^{(a)} \\ \mathbf{T}_{6,5}^{(b)} & \mathbf{T}_{6,6}^{(b)} + \mathbf{T}_{6,6}^{(c)} & \mathbf{T}_{6,7}^{(c)} & \mathbf{0} \\ \mathbf{0} & \mathbf{T}_{7,6}^{(c)} & \mathbf{T}_{7,7}^{(c)} + \mathbf{T}_{7,7}^{(d)} & \mathbf{T}_{7,8}^{(d)} \\ \mathbf{T}_{8,5}^{(a)} & \mathbf{0} & \mathbf{T}_{8,7}^{(d)} & \mathbf{T}_{8,8}^{(d)} + \mathbf{T}_{8,8}^{(a)} \end{bmatrix}. \quad (28)$$

Now that the matrices and vectors are defined, we can construct the linear system in Equation (3) and compute the merged data.

## B Full algorithms for 2D problems using ItI matrices

In this section, we describe the details for the two-dimensional version of our method which merges ItI matrices. In this version of the algorithm, the outgoing boundary data  $\mathbf{h}$  tabulates the outgoing impedance data due to the particular solution, and the incoming boundary data

$\mathbf{g}$  tabulates incoming impedance data due to the homogeneous solution.  $\mathbf{T}$  is an impedance-to-impedance matrix. To define the impedance data, we need to choose a value  $\eta \in \mathbb{R}_+$ . In the wave scattering context, we often choose  $\eta = k$  (Gillman et al., 2015).

As in the previous section, we use  $\mathbf{I}_{a \times a}$  to denote the identity matrix of shape  $a \times a$  and  $\mathbf{0}_d$  to denote a length- $d$  vector filled with 0's. We also use  $\mathbf{0}$  to denote a block filled with 0's, and assume the shape of this block can be inferred from the context.

### B.1 Local solve stage

As before, there are  $4p - 4$  Chebyshev points on the boundary and  $4q$  Gauss points on the boundary, and we must map between these two sets of discretization points. In particular, we need to precompute the following matrices:

- $\mathbf{P}$ , with shape  $4p - 4 \times 4q$ , is the operator mapping data sampled on the Gauss boundary points to data sampled on the  $4p - 4$  Chebyshev points located on the boundary of the leaf. This matrix is constructed using a barycentric Lagrange interpolation matrix mapping from Gauss to Chebyshev points on one side of the leaf with the final row deleted; this interpolation matrix is repeated for the other sides.
- $\mathbf{Q}$ , with shape  $4q \times 4p$ , is the operator mapping data sampled on the Chebyshev points located on the boundary of the leaf to the Gauss points on the boundary of the leaf. This matrix is block-diagonal with four copies of a barycentric Lagrange interpolation matrix mapping from Chebyshev to Gauss points on one side of the leaf. Note this matrix double-counts the Chebyshev points at the corners of the leaf.
- $\mathbf{N}$ , with shape  $4p \times p^2$ , is an operator mapping from interior solutions to outward-pointing normal derivative data evaluated on the boundary Chebyshev points. Note this operator counts each corner point twice. This matrix is formed by stacking relevant rows of Chebyshev spectral differentiation matrices.
- $\tilde{\mathbf{N}}$ , with shape  $4p - 4 \times p^2$ , is an operator mapping from interior solutions to outward-pointing normal data evaluated on the boundary Chebyshev points. Note this operator counts each corner point only once. This matrix is formed by stacking relevant rows of Chebyshev spectral differentiation matrices.
- $\mathbf{H}$ , with shape  $4p \times p^2$ , is an operator mapping from interior solutions to evaluations of outgoing impedance data on the Chebyshev boundary discretization points. This matrix is constructed by taking  $\mathbf{N}$  and subtracting  $i\eta \mathbf{I}_{p \times p}$  from the appropriate submatrices. Again, this matrix double-counts the Chebyshev discretization points at the corners of the leaf.
- $\mathbf{G}$ , with shape  $4p - 4 \times p^2$ , is an operator mapping from interior solutions to evaluations of incoming impedance data on the Chebyshev boundary discretization points. This matrix is constructed by taking  $\tilde{\mathbf{N}}$  and adding  $i\eta \mathbf{I}_{p-1 \times p-1}$  to the appropriate submatrices.

To work with  $\mathbf{L}^{(i)}$ , the discretization of the differential operator on leaf  $i$ , it is useful to identify  $I_i$  and  $I_e$ , the sets of discretization points corresponding to the  $(p - 2)^2$  interior and  $4p - 4$  exterior Chebyshev points, respectively. As in Section A, we solve the local problem by using these precomputed operators to enforce the differential operator on the interior discretization points and the boundary condition on the boundary discretization points. In this case, the boundary condition is an incoming impedance condition, also known as a Robin boundary condition. Now, we can fully describe the local solve stage in Algorithm 8.

---

**Algorithm 8:** 2D ItI local solve stage.

---

**Input:** Discretized differential operators  $\{\mathbf{L}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$ ; discretized source functions  $\{\mathbf{f}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$ ; precomputed interpolation and differentiation matrices  $\mathbf{P}, \mathbf{Q}, \mathbf{H}$ , and  $\mathbf{G}$ .

- 1 **for**  $i = 1, \dots, n_{\text{leaves}}$  **do**
- 2    $\mathbf{B}^{(i)} = \begin{bmatrix} \mathbf{G} \\ \mathbf{L}^{(i)}(I_i, :) \end{bmatrix}$
- 3   Invert  $\mathbf{B}^{(i)}$
- 4    $\mathbf{Y}^{(i)} = (\mathbf{B}^{(i)})^{-1}(:, I_e) \mathbf{P}$
- 5    $\mathbf{v}^{(i)} = (\mathbf{B}^{(i)})^{-1}(:, I_i) \mathbf{f}^{(i)}(I_i)$
- 6    $\mathbf{T}^{(i)} = \mathbf{Q} \mathbf{H} \mathbf{Y}^{(i)}$
- 7    $\mathbf{h}^{(i)} = \mathbf{Q} \mathbf{H} \mathbf{v}^{(i)}$

**Result:** Poincaré–Steklov matrices  $\{\mathbf{T}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$ ; outgoing boundary data  $\{\mathbf{h}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$ ; interior solution matrices  $\{\mathbf{Y}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$ ; leaf-level particular solutions  $\{\mathbf{v}^{(i)}\}_{i=1}^{n_{\text{leaves}}}$

---

## B.2 Merge stage

As in Section A, we are merging four nodes  $\Omega_a, \Omega_b, \Omega_c$ , and  $\Omega_d$  with boundary parts labeled  $1, 2, \dots, 8$ . See Figure 12 for a diagram of the different boundary parts.

To implement the merge stage, we use sets of constraints to solve for a mapping from given  $\mathbf{g}_{\text{ext}}$  to unknown  $\mathbf{g}_{\text{int}}$ . These are vectors tabulating incoming impedance data due to the homogeneous solution along the boundary parts. Because  $\mathbf{g}_{\text{int}}$  tabulates impedance data, we must represent the incoming data with respect to neighboring nodes separately. For example, we must represent  $\mathbf{g}_5^{(a)}$ , the data along boundary element 5 incoming to node  $a$ , separately from  $\mathbf{g}_5^{(b)}$ . To that end, we use  $\mathbf{g}_{\text{int}} = [\mathbf{g}_5^{(a)}, \mathbf{g}_8^{(a)}, \mathbf{g}_6^{(c)}, \mathbf{g}_7^{(c)}, \mathbf{g}_5^{(b)}, \mathbf{g}_6^{(b)}, \mathbf{g}_7^{(d)}, \mathbf{g}_8^{(d)}]^\top$ . The ordering of these boundary elements is chosen specifically to reduce the computation during the merge, which will become apparent later. Again, we use  $n_{\text{side}}$  to denote the number of discretization points along each side of the nodes being merged, so  $\mathbf{g}_{\text{int}}$  has length  $8n_{\text{side}}$ .

The first set of constraints specify that the solution to the PDE is continuous:

$$\mathbf{u}_{\text{ext}} = \mathbf{A} \mathbf{g}_{\text{ext}} + \mathbf{B} \mathbf{g}_{\text{int}} + \mathbf{h}_{\text{ext}}^{(\text{child})}. \quad (29)$$

In this equation  $\mathbf{u}_{\text{ext}}$  is interpreted as the outgoing impedance data of the solution to the PDE restricted to the merged nodes with boundary data specified by  $\mathbf{g}_{\text{ext}}$ . In this set of constraints, we define:

$$\mathbf{h}_{\text{ext}}^{(\text{child})} = \begin{bmatrix} \mathbf{h}_1^{(a)} \\ \mathbf{h}_2^{(b)} \\ \mathbf{h}_3^{(c)} \\ \mathbf{h}_4^{(d)} \end{bmatrix}, \quad (30)$$

$$\mathbf{A} = \begin{bmatrix} \mathbf{T}_{1,1}^{(a)} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{T}_{2,2}^{(b)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{T}_{3,3}^{(c)} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{T}_{4,4}^{(d)} \end{bmatrix}, \quad (31)$$



$$B = \begin{bmatrix} T_{1,5}^{(a)} & T_{1,8}^{(a)} & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & T_{2,5}^{(b)} & T_{2,6}^{(b)} & 0 & 0 \\ 0 & 0 & T_{3,6}^{(c)} & T_{3,7}^{(c)} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & T_{4,7}^{(d)} & T_{4,8}^{(d)} \end{bmatrix}. \quad (32)$$

A second set of constraints specifies that the outgoing total solution's impedance data from one node must be opposite to the incoming homogeneous solution's impedance data for the neighboring node. For example, along merge interface 5, we enforce this constraint:

$$\mathbf{0}_{n_{\text{side}}} = \mathbf{u}_5^{(b)} + \mathbf{g}_5^{(a)}. \quad (33)$$

In this equation,  $\mathbf{u}_5^{(b)}$  is the outgoing impedance data due to the total solution of the PDE restricted to the merged nodes with boundary condition  $\mathbf{g}_{\text{ext}}$ . The normal derivative is oriented relative to node  $b$ . We can expand  $\mathbf{u}_5^{(b)}$  to find:

$$\mathbf{0}_{n_{\text{side}}} = \mathbf{g}_5^{(a)} + T_{5,5}^{(a)} \mathbf{g}_5^{(a)} + T_{5,8}^{(a)} \mathbf{g}_8^{(a)} + T_{5,5}^{(b)} \mathbf{g}_5^{(b)} + T_{5,6}^{(b)} \mathbf{g}_6^{(b)} + T_{5,5}^{(a)} \mathbf{g}_5^{(a)} + T_{5,1}^{(a)} \mathbf{g}_1^{(a)} + \mathbf{h}_5^{(a)}. \quad (34)$$

Similar equalities hold in each direction along each merge interface. We can expand to form a second system of constraints:

$$-\mathbf{h}_{\text{int}}^{(\text{child})} = \mathbf{C} \mathbf{g}_{\text{ext}} + \mathbf{D} \mathbf{g}_{\text{int}}, \quad (35)$$

where we define

$$\mathbf{h}_{\text{int}}^{(\text{child})} = \begin{bmatrix} \mathbf{h}_5^{(b)} \\ \mathbf{h}_8^{(d)} \\ \mathbf{h}_6^{(b)} \\ \mathbf{h}_7^{(d)} \\ \mathbf{h}_5^{(a)} \\ \mathbf{h}_6^{(c)} \\ \mathbf{h}_7^{(c)} \\ \mathbf{h}_8^{(a)} \end{bmatrix}, \quad (36)$$

$$\mathbf{C} = \begin{bmatrix} 0 & T_{5,2}^{(b)} & 0 & 0 \\ 0 & 0 & 0 & T_{8,4}^{(d)} \\ 0 & T_{6,2}^{(b)} & 0 & 0 \\ 0 & 0 & 0 & T_{7,4}^{(d)} \\ T_{5,1}^{(a)} & 0 & 0 & 0 \\ 0 & 0 & T_{6,3}^{(c)} & 0 \\ 0 & 0 & T_{7,3}^{(c)} & 0 \\ T_{8,1}^{(a)} & 0 & 0 & 0 \end{bmatrix}, \quad (37)$$

$$\mathbf{D} = \mathbf{I}_{8n_{\text{side}} \times 8n_{\text{side}}} + \begin{bmatrix} 0 & 0 & 0 & 0 & T_{5,5}^{(b)} & T_{5,6}^{(b)} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & T_{8,7}^{(d)} & T_{8,8}^{(d)} \\ 0 & 0 & 0 & 0 & T_{6,5}^{(b)} & T_{6,6}^{(b)} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & T_{7,7}^{(d)} & T_{7,8}^{(d)} \\ T_{5,5}^{(a)} & T_{5,8}^{(a)} & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & T_{6,6}^{(c)} & T_{6,7}^{(c)} & 0 & 0 & 0 & 0 \\ 0 & 0 & T_{7,6}^{(c)} & T_{7,7}^{(c)} & 0 & 0 & 0 & 0 \\ T_{8,5}^{(a)} & T_{8,8}^{(a)} & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}. \quad (38)$$

$D$  has a special structure which allows us to efficiently compute  $D^{-1}$  via Schur complement methods. Note that we can re-write Equation (38) as a block matrix with  $2 \times 2$  blocks:

$$D = \begin{bmatrix} I_{4n_{\text{side}} \times 4n_{\text{side}}} & D_{12} \\ D_{21} & I_{4n_{\text{side}} \times 4n_{\text{side}}} \end{bmatrix}. \quad (39)$$

This structure allows us to construct  $W = I_{4n_{\text{side}} \times 4n_{\text{side}}} - D_{12}D_{21}$ , the Schur complement of the lower-right  $I_{4n_{\text{side}} \times 4n_{\text{side}}}$  block in  $D$ ; this is the only matrix we need to invert to compute  $D^{-1}$ :

$$D^{-1} = \begin{bmatrix} W^{-1} & -W^{-1}D_{12} \\ -D_{21}W^{-1} & I_{4n_{\text{side}} \times 4n_{\text{side}}} + D_{21}W^{-1}D_{12} \end{bmatrix}. \quad (40)$$

We use Equation (40) to compute  $D^{-1}$  and then construct the outputs of the merge stage using Equations (4) and (5).

## C Full algorithms for 3D problems using DtN matrices with a uniform discretization

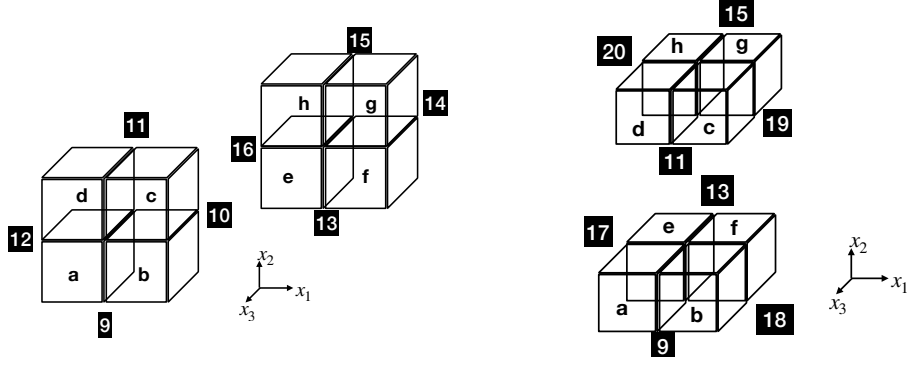
In this section, we describe the details for the three-dimensional version of our method which merges DtN matrices. In this version of the algorithm, the outgoing boundary data  $\mathbf{h}$  tabulates the outward-pointing normal derivative of the particular solution, and the incoming boundary data  $\mathbf{g}$  tabulates the homogenous solution values restricted to patch boundaries.  $\mathbf{T}$  is a Dirichlet-to-Neumann matrix.

As in the previous section, we use  $I_{a \times a}$  to denote the identity matrix of shape  $a \times a$  and  $\mathbf{0}_d$  to denote a length- $d$  vector filled with 0's. We also use  $\mathbf{0}$  to denote a matrix block filled with 0's, and assume the shape of this block can be inferred from its context.

### C.1 Local solve stage

Recall from Section 3.1 that we use a tensor product of order- $p$  Chebyshev–Lobatto points to discretize the interior of each leaf. This results in a grid with  $p^3$  discretization points, and  $p^3 - (p - 2)^3$  of these points lie on the boundary of the leaf. We use order- $q$  Gauss–Legendre points to discretize each side of the leaf's boundary, so there are  $6q^2$  boundary points in total. Thus, to translate the information between the interior and boundary of each leaf, we need to compute spectral differentiation matrices and matrices interpolating between the  $p^2$  Chebyshev and  $q^2$  Gauss points on each face of the leaf. In particular, we need to precompute the following matrices:

- $\mathbf{P}$ , with shape  $p^3 - (p - 2)^3 \times 6q^2$ , is the operator mapping data sampled on the Gauss boundary points to data sampled on the Chebyshev points located on the boundary of the leaf. This matrix is constructed using a barycentric Lagrange interpolation matrix mapping from Gauss to Chebyshev points on one face of the leaf; this interpolation matrix is repeated for the other sides. Rows corresponding to the Chebyshev points on the corners (edges) of the leaf average the contribution from the three (two) adjoining panels.
- $\mathbf{Q}$  with shape  $6q^2 \times p^3$ , performs spectral differentiation on the  $p^3$  Chebyshev points followed by interpolation to the Gauss boundary points. This matrix is formed by stacking the relevant rows of Chebyshev spectral differentiation matrices to form an operator which evaluates normal derivatives on the boundary Chebyshev points, and then composing this differentiation operator with a matrix formed from barycentric Lagrange interpolation



**Figure 13:** Visualizing boundary elements 9 through 20 for three-dimensional merges.

matrix blocks. These interpolation matrices each map from one Chebyshev face to one Gauss face.

To work with  $\mathbf{L}^{(i)}$ , the discretization of the differential operator on leaf  $i$ , it is useful to identify  $I_i$  and  $I_e$ , the sets of discretization points corresponding to the  $(p-2)^3$  interior and  $p^3 - (p-2)^3$  exterior Chebyshev points, respectively. Once the precomputed operators and index sets are correctly specified, Algorithm 7 can be re-used for the three-dimensional case.

## C.2 Merge stage

In the 3D merge stage, we are merging eight nodes  $\Omega_a, \Omega_b, \Omega_c, \Omega_d, \Omega_e, \Omega_f, \Omega_g$ , and  $\Omega_h$ , which have *exterior* and *interior* discretization points. We label the exterior boundary sections  $1, 2, \dots, 8$ , and we label the interior boundary sections  $9, 10, \dots, 20$ . See Figure 13 for a diagram of the different boundary parts. Because the merge stage operates completely on data discretized using Gauss–Legendre panels, there are no discretization points at the corners or edges of nodes. This means each discretization point belongs to exactly one part of the boundary. During this stage of the algorithm, we will be indexing rows and columns of the Dirichlet-to-Neumann matrices according to boundary sections  $1, \dots, 20$ . For example, we use  $\mathbf{T}_{1,9}^{(a)}$  to indicate the submatrix of node  $a$ ’s DtN matrix which maps from boundary section 9 to boundary section 1. Suppose that each node has  $n_{\text{side}}$  discretization points on each face. Then  $\mathbf{T}_{1,9}^{(a)}$  will have shape  $3n_{\text{side}} \times n_{\text{side}}$ .

Just as in the two-dimensional case, we use sets of constraints to solve for a mapping from given  $\mathbf{g}_{\text{ext}}$  to unknown  $\mathbf{g}_{\text{int}}$ , vectors tabulating the homogeneous solution along the boundary parts. First are constraints specifying that the solution to the PDE is continuous:

$$\mathbf{u}_{\text{ext}} = \mathbf{A}\mathbf{g}_{\text{ext}} + \mathbf{B}\mathbf{g}_{\text{int}} + \mathbf{h}_{\text{ext}}^{(\text{child})}. \quad (41)$$

In this equation  $\mathbf{u}_{\text{ext}}$  is interpreted as the outward-pointing normal derivative of the solution to the PDE restricted to the merged nodes with boundary data specified by  $\mathbf{g}_{\text{ext}}$ . In this set of

constraints, we define:

$$\mathbf{h}_{\text{ext}}^{(\text{child})} = \begin{bmatrix} \mathbf{h}_1^{(a)} \\ \mathbf{h}_2^{(b)} \\ \mathbf{h}_3^{(c)} \\ \mathbf{h}_4^{(d)} \\ \mathbf{h}_5^{(e)} \\ \mathbf{h}_6^{(f)} \\ \mathbf{h}_7^{(g)} \\ \mathbf{h}_8^{(h)} \end{bmatrix}, \quad (42)$$

$$\mathbf{A} = \begin{bmatrix} \mathbf{T}_{1,1}^{(a)} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{T}_{2,2}^{(b)} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{T}_{3,3}^{(c)} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{T}_{4,4}^{(d)} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{T}_{5,5}^{(e)} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{T}_{6,6}^{(f)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{T}_{7,7}^{(g)} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{T}_{8,8}^{(h)} \end{bmatrix}, \quad (43)$$

$$\mathbf{B} = \begin{bmatrix} \mathbf{T}_{1,9}^{(a)} & \mathbf{0} & \mathbf{0} & \mathbf{T}_{1,12}^{(a)} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{T}_{1,17}^{(a)} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{T}_{2,9}^{(b)} & \mathbf{T}_{2,10}^{(b)} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{T}_{2,18}^{(b)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{T}_{3,10}^{(c)} & \mathbf{T}_{3,11}^{(c)} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{T}_{3,19}^{(c)} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{T}_{4,11}^{(d)} & \mathbf{T}_{4,12}^{(d)} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{T}_{4,20}^{(d)} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{T}_{5,13}^{(e)} & \mathbf{0} & \mathbf{0} & \mathbf{T}_{5,16}^{(e)} & \mathbf{T}_{5,17}^{(e)} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{T}_{6,13}^{(f)} & \mathbf{T}_{6,14}^{(f)} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{T}_{6,18}^{(f)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{T}_{7,14}^{(g)} & \mathbf{T}_{7,15}^{(g)} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{T}_{7,19}^{(g)} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{T}_{8,15}^{(h)} & \mathbf{T}_{8,16}^{(h)} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{T}_{8,20}^{(h)} \end{bmatrix}. \quad (44)$$

As in the two-dimensional case, we enforce constraints that ensure the normal derivatives from neighboring nodes sum to zero, which gives us a system of constraints:

$$\mathbf{0}_{12n_{\text{side}}} = \mathbf{C}\mathbf{g}_{\text{ext}} + \mathbf{D}\mathbf{g}_{\text{int}} + \mathbf{h}_{\text{int}}^{(\text{child})}, \quad (45)$$

where we define

$$\mathbf{h}_{\text{int}}^{(\text{child})} = \begin{bmatrix} \mathbf{h}_9^{(a)} + \mathbf{h}_9^{(b)} \\ \mathbf{h}_{10}^{(b)} + \mathbf{h}_{10}^{(c)} \\ \mathbf{h}_{11}^{(c)} + \mathbf{h}_{11}^{(d)} \\ \mathbf{h}_{12}^{(d)} + \mathbf{h}_{12}^{(a)} \\ \mathbf{h}_{13}^{(e)} + \mathbf{h}_{13}^{(f)} \\ \mathbf{h}_{14}^{(f)} + \mathbf{h}_{14}^{(g)} \\ \mathbf{h}_{15}^{(g)} + \mathbf{h}_{15}^{(h)} \\ \mathbf{h}_{16}^{(h)} + \mathbf{h}_{16}^{(e)} \\ \mathbf{h}_{17}^{(a)} + \mathbf{h}_{17}^{(e)} \\ \mathbf{h}_{18}^{(b)} + \mathbf{h}_{18}^{(f)} \\ \mathbf{h}_{19}^{(c)} + \mathbf{h}_{19}^{(g)} \\ \mathbf{h}_{20}^{(d)} + \mathbf{h}_{20}^{(h)} \end{bmatrix}, \quad (46)$$

$$\mathbf{C} = \begin{bmatrix} \mathbf{T}_{9,1}^{(a)} & \mathbf{T}_{9,2}^{(b)} & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & \mathbf{T}_{10,2}^{(b)} & \mathbf{T}_{10,3}^{(c)} & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & \mathbf{T}_{11,3}^{(c)} & \mathbf{T}_{11,4}^{(d)} & 0 & 0 & 0 & 0 \\ \mathbf{T}_{12,1}^{(a)} & 0 & 0 & \mathbf{T}_{12,4}^{(d)} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & \mathbf{T}_{13,5}^{(e)} & \mathbf{T}_{13,6}^{(f)} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & \mathbf{T}_{14,6}^{(f)} & \mathbf{T}_{14,7}^{(g)} & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & \mathbf{T}_{15,7}^{(g)} & \mathbf{T}_{15,8}^{(h)} \\ 0 & 0 & 0 & 0 & \mathbf{T}_{16,5}^{(e)} & 0 & 0 & \mathbf{T}_{16,8}^{(h)} \\ \mathbf{T}_{17,1}^{(a)} & 0 & 0 & 0 & \mathbf{T}_{17,5}^{(e)} & 0 & 0 & 0 \\ 0 & \mathbf{T}_{18,2}^{(b)} & 0 & 0 & 0 & \mathbf{T}_{18,6}^{(f)} & 0 & 0 \\ 0 & 0 & \mathbf{T}_{19,3}^{(c)} & 0 & 0 & 0 & \mathbf{T}_{19,7}^{(g)} & 0 \\ 0 & 0 & 0 & \mathbf{T}_{20,4}^{(d)} & 0 & 0 & 0 & \mathbf{T}_{20,8}^{(h)} \end{bmatrix}, \quad (47)$$

$$\mathbf{D} = \begin{bmatrix} \mathbf{T}_{9,9}^{(a)} + \mathbf{T}_{9,9}^{(b)} & \mathbf{T}_{10,9}^{(b)} & 0 & \mathbf{T}_{10,12}^{(a)} & 0 & 0 & 0 & 0 & \mathbf{T}_{10,17}^{(a)} & \mathbf{T}_{10,18}^{(b)} & 0 & 0 \\ \mathbf{T}_{10,9}^{(b)} & \mathbf{T}_{10,10}^{(b)} + \mathbf{T}_{10,10}^{(c)} & \mathbf{T}_{10,11}^{(c)} & 0 & 0 & 0 & 0 & 0 & 0 & \mathbf{T}_{10,18}^{(b)} & \mathbf{T}_{10,19}^{(c)} & 0 \\ 0 & \mathbf{T}_{11,10}^{(c)} & \mathbf{T}_{11,11}^{(c)} + \mathbf{T}_{11,11}^{(d)} & \mathbf{T}_{11,12}^{(d)} & 0 & 0 & 0 & 0 & 0 & \mathbf{T}_{11,19}^{(c)} & \mathbf{T}_{11,20}^{(d)} & \mathbf{T}_{12,20}^{(d)} \\ \mathbf{T}_{12,9}^{(d)} & 0 & \mathbf{T}_{12,11}^{(d)} & \mathbf{T}_{12,12}^{(d)} + \mathbf{T}_{12,12}^{(a)} & 0 & 0 & 0 & \mathbf{T}_{12,17}^{(c)} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & \mathbf{T}_{13,13}^{(c)} + \mathbf{T}_{13,13}^{(f)} & \mathbf{T}_{13,14}^{(f)} & 0 & \mathbf{T}_{13,16}^{(c)} & \mathbf{T}_{13,17}^{(c)} & \mathbf{T}_{13,18}^{(f)} & 0 & 0 \\ 0 & 0 & 0 & 0 & \mathbf{T}_{14,13}^{(f)} & \mathbf{T}_{14,14}^{(g)} + \mathbf{T}_{14,14}^{(g)} & \mathbf{T}_{14,15}^{(g)} & 0 & 0 & \mathbf{T}_{14,18}^{(f)} & \mathbf{T}_{14,19}^{(g)} & 0 \\ 0 & 0 & 0 & 0 & 0 & \mathbf{T}_{15,14}^{(g)} & \mathbf{T}_{15,15}^{(g)} + \mathbf{T}_{15,15}^{(h)} & \mathbf{T}_{15,16}^{(h)} & \mathbf{T}_{15,17}^{(c)} & \mathbf{T}_{15,18}^{(g)} & \mathbf{T}_{15,19}^{(h)} & \mathbf{T}_{16,20}^{(h)} \\ 0 & 0 & 0 & 0 & \mathbf{T}_{16,13}^{(c)} & 0 & \mathbf{T}_{16,15}^{(h)} & \mathbf{T}_{16,16}^{(h)} + \mathbf{T}_{16,16}^{(c)} & \mathbf{T}_{16,17}^{(c)} & 0 & 0 & 0 \\ \mathbf{T}_{17,9}^{(c)} & 0 & 0 & \mathbf{T}_{17,12}^{(a)} & \mathbf{T}_{17,13}^{(c)} & 0 & 0 & \mathbf{T}_{17,16}^{(c)} & \mathbf{T}_{17,17}^{(c)} + \mathbf{T}_{17,17}^{(c)} & 0 & 0 & 0 \\ \mathbf{T}_{18,9}^{(b)} & \mathbf{T}_{18,10}^{(b)} & 0 & 0 & \mathbf{T}_{18,13}^{(f)} & \mathbf{T}_{18,14}^{(f)} & 0 & 0 & 0 & \mathbf{T}_{18,18}^{(b)} + \mathbf{T}_{18,18}^{(f)} & 0 & 0 \\ 0 & \mathbf{T}_{19,10}^{(c)} & \mathbf{T}_{19,11}^{(c)} & 0 & 0 & \mathbf{T}_{19,14}^{(g)} & \mathbf{T}_{19,15}^{(g)} & 0 & 0 & 0 & \mathbf{T}_{19,19}^{(c)} + \mathbf{T}_{19,19}^{(g)} & 0 \\ 0 & 0 & \mathbf{T}_{20,11}^{(d)} & \mathbf{T}_{20,12}^{(d)} & 0 & 0 & \mathbf{T}_{20,15}^{(h)} & \mathbf{T}_{20,16}^{(h)} & 0 & 0 & 0 & \mathbf{T}_{20,20}^{(d)} + \mathbf{T}_{20,20}^{(h)} \end{bmatrix}. \quad (48)$$

Now that the matrices and vectors are defined, we can construct the linear system in Equation (3) and compute the merged data.

## D Details for the inverse scattering experiment

We consider a forward model which composes  $\mathcal{F}$ , which maps a scattering potential to evaluations of a scattered wave, and  $\mathcal{B}$ , which maps coefficients of a sine series to scattering potentials:

$$(\mathcal{F} \circ \mathcal{B})(\theta)_j = u_\theta(x^{(j)}); \quad u_\theta \text{ solves Equation 9 with } q = q_\theta; \quad (49)$$

$$q_\theta(x) = \mathcal{B}(\theta) = \sum_{b_j \in B_\gamma} \theta_j b_j(x); \quad (50)$$

where the basis  $B_\gamma$  is specified by Equation (12). In this experiment, we use a domain  $\Omega = [-1, 1]$ , and we use 100 evaluation points equispaced on a ring of radius 5:

$$x^{(j)} = \left( 5 \sin \left( \frac{2\pi j}{100} \right), 5 \cos \left( \frac{2\pi j}{100} \right) \right), \quad j = 0, \dots, 99. \quad (51)$$

To solve an inverse problem, we are interested in computing the Fréchet derivative of  $\mathcal{F} \circ \mathcal{B}$  centered at  $\theta$ ; we call this object  $J[\theta]$ . By the chain rule, we can decompose this Fréchet derivative

$$J[\theta] = J_{\mathcal{F}}[q_\theta] \circ J_{\mathcal{B}}[\theta].$$

The basis transformation  $\mathcal{B}$  is linear, so the action of  $J_{\mathcal{B}}[\theta]$  can be computed with standard sine and adjoint sine transforms. The following section will describe how we compute the action of  $J_{\mathcal{F}}[q_\theta]$ .

### D.1 Defining the action of the Fréchet derivative

Borges et al. (2017) describe the Fréchet derivative  $J_{\mathcal{F}}[q_\theta]$  the action  $J_{\mathcal{F}}[q_\theta]v$  and  $J_{\mathcal{F}}[q_\theta]^*f$ . The action of the derivative and its adjoint can be described by the solution of elliptic partial differential equations. We re-state these results in this section.

**Theorem D.1 (Theorem 3.1 of Borges et al. (2017))** *Let  $u_\theta$  solve Equation (9) with scattering potential  $q = q_\theta$ . Let  $w$  solve*

$$\begin{cases} \Delta w(x) + k^2(1 + q(x))w(x) = k^2 v(x) \left( u_\theta(x) + e^{ik\langle \hat{s}, x \rangle} \right), & x \in [-1, 1]^2, \\ \sqrt{r} \left( \frac{\partial w}{\partial r} - ikw \right) \rightarrow 0, & r = \|x\|_2 \rightarrow \infty. \end{cases} \quad (52)$$

*Then*

$$(J_{\mathcal{F}}[q_\theta]v)_j = w(x^{(j)}).$$

**Theorem D.2 (Theorem 3.2 of Borges et al. (2017))** *Let  $u_\theta$  solve Equation (9) with scattering potential  $q = q_\theta$ . Let  $f$  denote a singular charge distribution supported on the evaluation points  $\{x^{(j)}\}_{j=0, \dots, 99}$  viewed as a generalized function in  $\mathbb{R}^2$ . Let  $w$  solve*

$$\begin{cases} \Delta w(x) + k^2(1 + q(x))w(x) = k^2 f(x), & x \in \mathbb{R}^2, \\ \sqrt{r} \left( \frac{\partial w}{\partial r} + ikw \right) \rightarrow 0, & r = \|x\|_2 \rightarrow \infty. \end{cases} \quad (53)$$

*Then*

$$(J_{\mathcal{F}}[q_\theta]^*f)(x) = w(x) \overline{u_\theta(x)} + e^{ik\langle \hat{s}, x \rangle}.$$

### D.2 Experiment details

In this section, we describe the experimental setting used to generate Figure 9a. In these experiments, we use  $k = 20$  and basis  $B_\gamma$  with  $\gamma = 25$ , which gives us  $N_\theta = 465$  basis coefficients. We compute a coefficient vector  $\theta$  by projecting the scattering potential Equation (10) onto  $B_\gamma$ .

To measure the accuracy of the outputs of JAX Jacobian-vector products, we generate a random coefficient vector  $\delta$  distributed i.i.d. Gaussian and compute  $J[\theta]\delta$ . We compute this Jacobian-vector product for a range of discretization sizes, holding  $p = 16$  constant and varying  $L = 1, \dots, 5$ . We compare this with the action of the Fréchet derivative, which is computed by



Subtree depth	$L$	$N$	Runtime (s)	% of Peak FLOPS
5	8	16,777,216	5.86	6.68%
5	9	67,108,864	24.50	10.86%
6	8	16,777,216	4.35	10.58%
6	9	67,108,864	18.85	15.59%
7	8	16,777,216	4.02	14.86%
7	9	67,108,864	17.43	20.01%

**Table 3:** Evaluating the effect of the subtree height parameter.

first computing  $J_{\mathcal{B}}[\theta]\delta$  and then applying  $J_{\mathcal{F}}[q_{\theta}]$  from Theorem D.1, which is computed using parameters  $p = 16$  and  $L = 5$ . We measure the relative  $\ell_{\infty}$  error between the outputs of JAX Jacobian-vector products and the the action of the Fréchet derivative.

To measure the accuracy of the outputs of JAX vector-Jacobian products, we generate a random perturbation  $f \in \mathbb{C}^{100}$ . The real and imaginary parts of each component of  $f$  are distributed i.i.d. Gaussian. We compute this vector-Jacobian product for a range of discretization sizes, again holding  $p = 16$  constant and varying  $L = 1, \dots, 5$ . We compare this with the action of the Fréchet derivative which is computed by first evaluating  $J_{\mathcal{F}}[q_{\theta}]^*f$  via Theorem D.2, which is computed using parameters  $p = 16$  and  $L = 5$ , and then applying  $J_{\mathcal{B}}[\theta]^*$ . We measure the relative  $\ell_{\infty}$  error between the outputs of JAX vector-Jacobian products and the action of the Fréchet derivative.

We note the choice of evaluation points exterior to the computational domain (Equation (51)) is not necessary for the convergence of automatic differentiation; we choose these evaluation points to ensure numerical stability when computing  $J_{\mathcal{F}}[q_{\theta}]^*f$ . In preliminary experiments, we observed accurate vector-Jacobian products when the evaluation points were located at the HPS discretization points. In this setting, the inputs of the vector-Jacobian product routine must be scaled by the appropriate quadrature weights.

## E Additional timing results

In this appendix, we extend the results shown in Figure 3 to include our subtree recomputation method evaluated with different subtree depths. As a heuristic, we choose to use the subtree recomputation depth to be the maximum depth where all computations for Algorithms 1 and 2 can fit on a single GPU. In Figure 3, we consider the DtN version of the method with polynomial order  $p = 16$ ; this results in subtree depth 7. In Table 3, we show the results of choosing different subtree depth parameters. We measure the runtime for large problem sizes  $L = 8$  and 9; recomputation is necessary for these problem sizes.

As we decrease the subtree depth, we see runtimes increase, as more transfers between the GPU and host RAM are required.